

On High-dimensional and Low-rank Tensor Bandits

Chengshuai Shi, Cong Shen, and Nicholas D. Sidiropoulos
 University of Virginia
 Charlottesville, VA 22904, USA
 {cs7ync, cong, nikos}@virginia.edu

Abstract—Most existing studies on linear bandits focus on a one-dimensional characterization of the overall system. While being representative, this formulation may fail to model applications with high-dimensional but favorable structures, such as the low-rank tensor representation for recommender systems. To address this limitation, this work studies a general tensor bandits model, where actions and system parameters are represented by tensors as opposed to vectors, and we particularly focus on the case that the unknown system tensor is low-rank. A novel bandit algorithm, coined TOFU (Tensor Optimism in the Face of Uncertainty), is developed. TOFU first leverages flexible tensor regression techniques to estimate low-dimensional subspaces associated with the system tensor. These estimates are then utilized to convert the original problem to a new one with norm constraints on its system parameters. Lastly, a norm-constrained bandit subroutine is adopted by TOFU, which utilizes these constraints to avoid exploring the entire high-dimensional parameter space. Theoretical analyses show that TOFU improves the best-known regret upper bound by a multiplicative factor that grows exponentially in the system order. A novel performance lower bound is also established, which further corroborates the efficiency of TOFU.

I. INTRODUCTION

The multi-armed bandits (MAB) framework [1], [2] has attracted growing interest in recent years as it can characterize a broad range of applications requiring sequential decision-making. An active research area in MAB is linear bandits [3], [4], where the actions are characterized by feature vectors. While being representative, this one-dimensional (i.e., vectorized) formulation may fail to capture practical applications with high-dimensional but favorable structures. We use the recommender system model to illustrate this limitation. An online shopping platform needs an effective advertising mechanism for its products. However, instead of only deciding which item to promote (as typically considered in standard linear bandits studies), the marketer also needs to consider many other factors. For example, the marketer may plan where to place to promotion (e.g., on the sidebar or as a pop-up) and how to highlight the promotion (e.g., emphasizing the discounts or the product quality). The overall strategy with all these factors will determine the effectiveness of this promotion.

Traditional recommendation strategies often leverage tensor formulations to capture the joint decisions concerning many associated factors [5]–[7]. However, as mentioned, existing

The work of CSs was supported in part by the US National Science Foundation (NSF) under awards ECCS-2029978, ECCS-2143559, CNS-2002902, the Virginia Commonwealth Cyber Initiative (CCI) smart cities project, and the Bloomberg Data Science Ph.D. Fellowship. The work of NS was partially supported by NSF IIS-1908070.

TABLE I
RELATED WORKS AND REGRET COMPARISONS

Algorithm	Regret
Vectorized LinUCB [3]	$\tilde{O}(d^N \sqrt{T})$
Matricized ESTT/ESTS [9]	$\tilde{O}(d^{\lfloor \frac{N}{2} \rfloor} r^{\lfloor \frac{N}{2} \rfloor} \sqrt{T})$
Tensor Elim. [10]; modified to general actions	$\tilde{O}(d^{N-1} r \sqrt{T})$
TOFU (Corollary 1)	$\tilde{O}(d^2 r^{N-2} \sqrt{T})$
Lower bound (Theorem 2)	$\Omega(r^N \sqrt{T})$

The time horizon is T . The considered system tensor is order- N and of size (d, d, \dots, d) . It also has a multi-linear rank (r, r, \dots, r) , where $r \leq d$.

bandits strategies are largely restricted to vectorized systems. Although vectorizing multi-dimensional systems can preserve element-wise information, structural information is often lost. Especially, as recognized in [5]–[7], tensors formulated to characterize recommender systems often process the attractive property of *low-rankness* which, however, no longer exists in the vectorized systems and thus cannot be exploited.

In this work, we study a general problem of tensor bandits for online decision-making, which extends the standard one-dimensional setting of linear bandits to a multi-dimensional and multi-linear one. In particular, each action is represented by a tensor (as opposed to a vector), and the mean reward of playing an action is the inner product between its feature tensor and an unknown system tensor. Then, motivated by various practical problems, a low-rank assumption is imposed on the system tensor, and this work aims at leveraging the low-rank knowledge to facilitate bandit learning. The main contributions are summarized in the following.

- The studied tensor bandits framework is general in the sense that it does not have restrictions on the system dimension and the action structure, which contributes to the generalization of linear bandits and extends the applicability of the MAB study; see Appendix A for related works.

- A novel learning algorithm, TOFU (Tensor Optimism in the Face of Uncertainty), is proposed for the challenging problem of low-rank tensor bandits. TOFU adopts flexible designs of tensor regressions to estimate low-dimensional subspaces associated with the unknown system tensor. Then, these estimates are utilized to convert the original problem into a new one, where the low-rank property is transformed into the knowledge of norm constraints on the system parameters. TOFU finally adopts the LowOFUL subroutine [8] to incorporate these norm constraints in bandit learning to avoid exploring the entire high-dimensional parameter space.

- Theoretical analyses demonstrate the effectiveness and efficiency of TOFU with performance guarantees. In particular,

the regret of TOFU improves the best-known regret upper bound by a multiplicative factor of order $O((d/r)^{\lceil N/2 \rceil - 2})$, where N is the order of the considered system tensor, d is the length of its modes, and $r \leq d$ denotes its multi-linear rank. Note that this improvement becomes more significant in high-dimensional problems, i.e., growing exponentially w.r.t. N . A novel regret lower bound is further established, and TOFU is shown to be sub-optimal only up to a factor of $O((d/r)^2)$, which does not scale with N . The baselines and the main results are summarized in Table I.

II. PROBLEM FORMULATION

A. Preliminaries on Tensors

An order- N tensor $\mathcal{Y} \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_N}$ has $\prod_{n \in [N]} d_n$ elements and can be viewed as a hyper-rectangle with edges (referred to as modes) of lengths (d_1, d_2, \dots, d_N) (see [11], [12] for comprehensive reviews). The tensor elements are identified to by their indices along each mode, e.g., $\mathcal{Y}_{i_1, i_2, \dots, i_N}$ denotes the (i_1, i_2, \dots, i_N) -th element of \mathcal{Y} , while a block is denoted by the index set of its contained elements, e.g., the block $\mathcal{Y}_{I_1, I_2, \dots, I_N}$ represents the elements with indices $(i_1, i_2, \dots, i_N) \in I_1 \times I_2 \times \dots \times I_N$. Moreover, fibers are one-dimensional sections of a tensor (as rows and columns in a matrix); thus an order- N tensor has N types of fibers.

Tensor operations. The inner product between tensor \mathcal{Y} and a same-shape tensor $\mathcal{B} \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_N}$ is the sum of the products of their elements:

$$\langle \mathcal{B}, \mathcal{Y} \rangle = \sum_{i_1 \in [d_1]} \sum_{i_2 \in [d_2]} \dots \sum_{i_N \in [d_N]} \mathcal{B}_{i_1, i_2, \dots, i_N} \mathcal{Y}_{i_1, i_2, \dots, i_N}.$$

The Frobenius norm is then defined as $\|\mathcal{Y}\|_F := \sqrt{\langle \mathcal{Y}, \mathcal{Y} \rangle}$.

The mode- n (matrix) product $\mathcal{Y} \times_n B$ between tensor \mathcal{Y} and matrix $B \in \mathbb{R}^{d'_n \times d_n}$ outputs an order- N tensor of size $(d_1, \dots, d_{n-1}, d'_n, d_{n+1}, \dots, d_N)$ with elements:

$$(\mathcal{Y} \times_n B)_{i_1, \dots, i_{n-1}, i'_n, i_{n+1}, \dots, i_N} = \sum_{i_n \in [d_n]} B_{i'_n, i_n} \mathcal{Y}_{i_1, \dots, i_n, \dots, i_N}.$$

In addition, matricization is the process of reordering tensor elements into a matrix. The mode- n matricization of tensor \mathcal{Y} is denoted as $\mathcal{M}_n(\mathcal{Y})$, whose columns are mode- n fibers of tensor \mathcal{Y} and dimensions are $(d_n, \prod_{n' \in [N] \setminus \{n\}} d_{n'})$. Similarly, vectorization converts a tensor to a vector with all its elements, which is denoted as $\text{vec}(\mathcal{Y})$ for tensor \mathcal{Y} .

Tucker decomposition. Similarly to matrices, tensor decomposition is a useful tool to characterize the structure of tensors. In this work, we mainly focus on the Tucker decomposition illustrated as follows: for tensor \mathcal{Y} , with r_n denoting the rank of its mode- n matricization, i.e., $r_n = \text{rank}(\mathcal{M}_n(\mathcal{Y}))$, and U_n the corresponding left singular vectors of $\mathcal{M}_n(\mathcal{Y})$, there exists a core tensor $\mathcal{G} \in \mathbb{R}^{r_1 \times r_2 \times \dots \times r_N}$ such that

$$\mathcal{Y} = \mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 \dots \times_N U_N =: \mathcal{G} \times_{n \in [N]} U_n,$$

which can be denoted as $\mathcal{Y} = [[\mathcal{G}; U_1, \dots, U_N]]$, and the tuple (r_1, \dots, r_N) is called the multi-linear rank of tensor \mathcal{Y} .

Additional notations. Typically, lowercase characters (e.g., x) stand for scalars while vectors are denoted with bold lowercase

characters (e.g., \mathbf{x}). Capital characters (e.g., X) are used for matrices, and calligraphic capital characters (e.g., \mathcal{X}) for tensors. In addition, $\|\cdot\|_2$ denotes the Euclidean norm for vectors and the spectral norm for matrices; for a vector \mathbf{y} and a matrix Γ , we denote $\|\mathbf{y}\|_\Gamma := \sqrt{\mathbf{y}^\top \Gamma \mathbf{y}}$.

B. Tensor Bandits

This work considers the following multi-dimensional bandit problem. At each time step $t \in [T]$, the player has access to an action set $\mathbb{A}_t \subseteq \mathbb{R}^{d_1 \times d_2 \times \dots \times d_N}$, i.e., the elements are tensors of size (d_1, d_2, \dots, d_N) . She needs to select one action \mathcal{A}_t from the set \mathbb{A}_t , and this action would bring her a reward of

$$r_t = \langle \mathcal{A}_t, \mathcal{X} \rangle + \varepsilon_t, \quad (1)$$

where $\mathcal{X} \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_N}$ is an unknown tensor of system parameters and ε_t is an independent 1-sub-Gaussian noise. We further denote $\mu_{\mathcal{A}} := \langle \mathcal{A}, \mathcal{X} \rangle$ as the expected reward of action \mathcal{A} and, without loss of generality, assume that $\|\mathcal{X}\|_F \leq C$ for $C > 0$ and $\max\{\|\mathcal{A}\|_F : \mathcal{A} \in \cup_{t \in [T]} \mathbb{A}_t\} \leq 1$.

The agent's objective is to minimize her regret against the per-step optimal actions $\mathcal{A}_t^* := \arg \max_{\mathcal{A} \in \mathbb{A}_t} \langle \mathcal{A}, \mathcal{X} \rangle$ [1]:

$$R(T) := \sum_{t \in [T]} (\langle \mathcal{A}_t^*, \mathcal{X} \rangle - \langle \mathcal{A}_t, \mathcal{X} \rangle).$$

C. The Low-rank Structure

It is possible to view the above problem as a $\prod_{n \in [N]} d_n$ -dimensional linear bandits problem by vectorizing the action tensor \mathcal{A}_t and the system tensor \mathcal{X} , which can then be solved by known algorithms [3], [4]. However, the high-dimensional structures of this system are not preserved by vectorization. Especially, one of the most commonly observed structures in real-world applications (e.g., recommender systems [5]–[7] and healthcare [13]–[16]) is the low-rankness. We give the general multi-linear rank assumption of \mathcal{X} as follows.

Assumption 1. *The unknown system tensor \mathcal{X} has a multi-linear rank of (r_1, r_2, \dots, r_N) and can be decomposed as $\mathcal{X} = [[\mathcal{G}; U_1, U_2, \dots, U_N]]$.*

To simplify the notations, in the following, it is assumed that $d_1 = \dots = d_N = d$ while $r_1 = \dots = r_N = r$. In practice, the rank r is often much smaller than the mode length d , especially for very large d . Hence, the following problem is at the center of this work: *can bandit algorithms be designed to exploit the low-rank structure of the system tensor?* Especially, the key question is how much performance improvement we can achieve, compared with the naive regret of $\tilde{O}(d^N \sqrt{T})$ [3] that is obtained by directly vectorizing the actions and the system.

Note that the design and analysis can be extended to the general case of $d_1 \neq \dots \neq d_N$ and $r_1 \neq \dots \neq r_N$ with minor notation modifications. Also, without loss of generality, it is assumed that N is of order $O(1)$ (i.e., a constant) and $N \geq 3$.

III. THE TOFU ALGORITHM

The TOFU algorithm (presented in Alg. 1) has two phases: A and B. Phase A aims at estimating the unknown system tensor \mathcal{X} up to a certain precision, especially its low-dimensional subspaces. With this estimate, the original bandit problem can be

has a norm that scales with $\tilde{O}((\eta(T_1))^k)$ (see Lemma 2), where the notation $: r$ denotes the set $[r]$ while $r + 1$: represents the set $[r + 1 : d]$ (thus the above block denotes the $r^{N-k}(d-r)^k$ tensor elements with indices $(i_1, \dots, i_{N-k}, i_{N-k+1}, \dots, i_N) \in [r] \times \dots \times [r] \times [r+1, d] \times \dots \times [r+1, d]$). This property holds similarly for other symmetrical blocks. As $\eta(T_1)$ typically decays with T_1 (because the estimation quality should increase with more data samples), the norm of the above block will become smaller as the length of Phase A increases, which can be captured by a norm constraint that will be described later.

To ease the exposition, we refer to the above block and its symmetrical ones as blocks with k tails, meaning the indices of their elements have k modes in the interval $[r+1 : d]$ (i.e., the tail). An illustration of these blocks in an order-3 tensor is provided in Fig. 1. Furthermore, the number of tensor elements in blocks with less than k tails is denoted as

$$q(k) := \sum_{i=0}^{k-1} \binom{N}{i} r^{N-i} (d-r)^i, \quad (6)$$

which is an important quantity in later designs and analyses.

Remark 1. Compared with previous works on matrix and tensor bandits [8]–[10], [21], the essence of this work is the observation that norm constraints commonly exist for blocks with different numbers of tails. In particular, [10] directly extends [8], [21] and only leverages the norm constraint on the block with N tails. Instead, Section IV will illustrate that the norm constraints on blocks with at least three tails can be leveraged together under a suitable $\eta(T_1)$, which then leads to the obtained performance improvement.

Algorithm 1 TOFU

Input: T ; rank r ; dimension N and d ; tensor regression alg. TRalg ; length of Phase A T_1 ; confidence parameter δ ; tails ρ

- 1: Sample $\mathcal{A}_t \in \mathbb{A}_t$ following the arm selection rule required by $\text{TRalg}(\cdot)$ and observe reward r_t , for $t \in [T_1]$ ▷ *Phase A*
- 2: Estimate $\hat{\mathcal{X}} = [[\hat{\mathcal{G}}; \hat{U}_1, \dots, \hat{U}_N]]$ with TRalg using $\mathcal{D}_A = \{(\mathcal{A}_t, r_t) : t \in [T_1]\}$, i.e., $\hat{\mathcal{X}} \leftarrow \text{TRalg}(\mathcal{D}_A)$
- 3: Set $C_\perp, \lambda, \lambda_\perp$ as in Theorem 1 ▷ *Phase B*
- 4: Initialize $\Lambda(\rho) \leftarrow \text{diag}(\lambda, \dots, \lambda, \lambda_\perp, \dots, \lambda_\perp)$, where the first $q(\rho)$ elements are λ ; $\Psi_{T_1} \leftarrow \{\mathbf{y} \in \mathbb{R}^{d^N} : \|\mathbf{y}\|_2 \leq C\}$
- 5: **for** $t = T_1 + 1, \dots, T$ **do**
- 6: Set $\hat{\mathbb{B}}_t \leftarrow \{\hat{\mathcal{B}}_t = \mathcal{A}_t \times_{n \in [N]} [\hat{U}_n, \hat{U}_{n,\perp}]^\top : \mathcal{A}_t \in \mathbb{A}_t\}$
- 7: Get $\hat{\mathbf{b}}_t \leftarrow \arg \max_{\hat{\mathcal{B}}_t \in \text{vec}(\hat{\mathbb{B}}_t)} \max_{\mathbf{y} \in \Psi_{t-1}} \langle \hat{\mathbf{b}}_t, \mathbf{y} \rangle$
- 8: Pull arm \mathcal{A}_t corresponding to $\hat{\mathbf{b}}_t$ and obtain reward r_t
- 9: Update \hat{B}_t with rows $\{\mathbf{b}_\tau^\top : \tau \in (T_1, t]\}$
- 10: Update \mathbf{r}_t with elements $\{r_\tau : \tau \in (T_1, t]\}$
- 11: Update $V_t \leftarrow \Lambda(\rho) + \hat{B}_t^\top \hat{B}_t$ and $\hat{\mathbf{y}} \leftarrow V_t^{-1} \hat{B}_t^\top \mathbf{r}_t$
- 12: Update $\sqrt{\beta_t} \leftarrow \sqrt{\log(\frac{\det(V_t)}{\det(\Lambda(\rho))\delta^2})} + \sqrt{\lambda}C + \sqrt{\lambda_\perp}C_\perp$
- 13: Update $\Psi_t \leftarrow \{\hat{\mathbf{y}} \in \mathbb{R}^{d^N} : \|\hat{\mathbf{y}} - \hat{\mathbf{y}}\|_{V_t} \leq \sqrt{\beta_t}\}$
- 14: **end for**

C. Phase B: Solving the Norm-constrained Linear Bandits

As illustrated above, after the projection, norm constraints can be obtained on some blocks of tensor $\hat{\mathcal{Y}}$. For flexibility, we consider that Phase B aims to leverage such constraints

on blocks with at least ρ tails, which contain $d^N - q(\rho)$ elements. The parameter ρ is an input with its value in $[N]$ that requires careful designs to balance losses from two phases and will be specified in Sec. IV (e.g., selected as $\rho = 3$ in Corollary 1). Equivalently, there exist norm constraints on parts of the elements in the unknown vector

$$\hat{\mathbf{y}} := \text{vec}(\hat{\mathcal{Y}}) \in \mathbb{R}^{d^N}. \quad (7)$$

If the vectorization of $\hat{\mathcal{Y}}$ is performed first on the block with zero tail and then gradually on those with one and more tails (see Fig. 1 for an example), we can compactly express the norm constraint on blocks with at least ρ tails as

$$\|\hat{\mathbf{y}}_{q(\rho)+1:d^N}\|_2 \leq C_\perp, \quad (8)$$

where the parameter C_\perp will be specified later in Theorem 1. This condition can be interpreted as that there are approximately only $q(\rho)$ effective parameters in $\hat{\mathbf{y}}$ while the other parameters are nearly ignorable due to their constrained norm.

Then, a *norm-constrained linear bandits* problem with d^N parameters needs to be solved. In particular, the action set is $\Phi_t := \text{vec}(\hat{\mathbb{B}}_t) \subseteq \mathbb{R}^{d^N}$ at step t , where $\text{vec}(\hat{\mathbb{B}}_t) := \{\text{vec}(\hat{\mathcal{B}}) : \hat{\mathcal{B}} \in \hat{\mathbb{B}}_t\}$, and the expected reward for action $\hat{\mathbf{b}} \in \Phi_t$ is $\langle \hat{\mathbf{b}}, \hat{\mathbf{y}} \rangle$. Additionally, an important norm constraint on $\hat{\mathbf{y}}$, i.e., Eqn. (8), is available to the learner. Inspired by [22], the LowOFUL algorithm is designed in [8] to tackle such norm-constrained linear bandits. Especially, a weighted regularization is performed to estimate the system parameter: at time step t , the following estimate $\hat{\mathbf{y}}$ of $\hat{\mathbf{y}}$ is obtained as $\hat{\mathbf{y}} \leftarrow \arg \min_{\mathbf{y}} \|\hat{B}_t \mathbf{y} - \mathbf{r}_t\|_2^2 + \|\mathbf{y}\|_{\Lambda(\rho)}^2 = V_t^{-1} \hat{B}_t^\top \mathbf{r}_t$, where matrix $\hat{B}_t \in \mathbb{R}^{t \times d^N}$ is constructed with previous action vectors $\{\hat{\mathbf{b}}_\tau : \tau \in (T_1, t]\}$ as rows, vector $\mathbf{r}_t \in \mathbb{R}^t$ has elements $\{r_\tau : \tau \in (T_1, t]\}$, matrix $\Lambda(\rho) = \text{diag}(\lambda, \dots, \lambda, \lambda_\perp, \dots, \lambda_\perp)$ (with λ as the first $q(\rho)$ elements and λ_\perp as the others), and $V_t = \Lambda(\rho) + \hat{B}_t^\top \hat{B}_t$. Then, an OFU-style arm-selection subroutine is adopted (lines 5–14 of Alg. 1).

Remark 2. To better understand the projection performed in Eqns. (2) and (4), an ideal scenario is considered where the decomposition matrices (U_1, \dots, U_N) are *exactly* known. Then, the projected action $\hat{\mathcal{B}}$ and system parameter $\hat{\mathcal{Y}}$ both match their “exact” versions $\mathcal{B} = \mathcal{A} \times_{n=1}^N [U_n, U_{n,\perp}]^\top$ and $\mathcal{Y} = \mathcal{G} \times_{n \in [N]} ([U_n, U_{n,\perp}]^\top U_n) = \mathcal{G} \times_{n \in [N]} ([I_r, \mathbf{0}_{r \times (d-r)}]^\top)$. Although \mathcal{Y} has d^N elements, there are only r^N non-zero ones in \mathcal{G} . However, for $\hat{\mathcal{Y}}$ projected via the imperfect estimates $(\hat{U}_1, \dots, \hat{U}_N)$, we can only guarantee some blocks of elements have small norms instead of being exact nulls as in \mathcal{Y} .

IV. THEORETICAL ANALYSIS

In this section, we formally establish the theoretical guarantee of the TOFU algorithm. First, the following assumption is adopted on the minimum singular value of the matricized system tensor, which is commonly used in the study of matrix bandits [8], [9], [21] and tensor bandits [23].

Assumption 3. *It holds that $\min_{n \in [N]} \{\omega_{\min}(\mathcal{M}_n(\mathcal{X}))\} \geq \omega$ for some parameter $\omega > 0$, where $\omega_{\min}(\cdot)$ returns the minimum positive singular value of a matrix.*

Then, the following regret upper bound can be established.

Theorem 1. *Under Assumptions 1, 2 and 3, with probability at least $1 - \delta$, using $\rho \in [N]$ as input and $\lambda = C^{-2}$, $\lambda_{\perp} = \frac{T}{q(\rho) \log(1+T/\lambda)}$, $C_{\perp} = 2^{N/2} C(\eta(T_1))^{\rho} \omega^{-\rho}$, if T_1 is chosen such that $\eta(T_1) \leq \omega$, the regret of TOFU can be bounded as*

$$R(T) \leq \tilde{O}\left(CT_1 + d^{\rho-1} r^{N-\rho+1} \sqrt{T} + C(\eta(T_1))^{\rho} \omega^{-\rho} T\right).$$

It is worth noting that this theorem applies to any tensor regression technique satisfying Assumption 2 and any input ρ , which demonstrates the flexibility of TOFU. Furthermore, the above regret bound has three terms. The first term characterizes the dataset collection in Phase A. The second term represents the learning loss from the $q(\rho)$ major elements in Phase B. The third one is from the other $d^N - q(\rho)$ elements, which are nearly ignorable but still contribute to the regret.

According to function $\eta(T_1)$, parameters ρ and T_1 should be carefully selected such that the overall regret in Theorem 1 is minimized. Two specific tensor regression techniques from [23], [24] are considered to instantiate $\eta(T_1)$: the first one is established with the selected arms having sub-Gaussian elements, while the second selects random one-hot tensors as arms. To avoid complicated expressions, confidence parameters δ_1, δ_2 , threshold parameters ι_1, ι_2 and scale parameters c_1, c_2 are adopted in the following, whose values are independent of T_1 and can be found in the corresponding references.

Example 1 (Section 4.2 of [23]). *If $T_1 > \iota_1$, all elements of \mathcal{A}_t are i.i.d. drawn from $1/d^N$ -sub-Gaussian distributions, and ε_t is an independent standard Gaussian noise, with probability at least $1 - \delta_1$, an estimate $\hat{\mathcal{X}} = [[\hat{\mathcal{G}}; U_1, \dots, U_N]]$ can be obtained from the tensor regression algorithm proposed in [23] such that $\|\hat{\mathcal{X}} - \mathcal{X}\|_F^2 \leq c_1 d^N (dr + r^N)/T_1$.*

Example 2 (Corollary 2 of [24]). *If $T_1 > \iota_2$, \mathcal{A}_t is a random one-hot tensor, and ε_t is an independent 1-sub-Gaussian noise, with probability at least $1 - \delta_2$, an estimate $\hat{\mathcal{X}} = [[\hat{\mathcal{G}}; U_1, \dots, U_N]]$ can be obtained from the tensor regression algorithm proposed in [24] such that $\|\hat{\mathcal{X}} - \mathcal{X}\|_F^2 \leq c_2 d^N (dr + r^N)/T_1$.*

In these examples, it can be seen that Assumption 2 holds with a high probability for $\eta(T_1) = \tilde{O}(\sqrt{d^N (dr + r^N)/T_1})$. Then, Theorem 1 leads to the following corollary.

Corollary 1. *Under Assumptions 1 and 3, if the conditions in Example 1 (resp. Example 2) can be satisfied in Phase A, using the tensor regression algorithm from [23] (resp. [24]) as $\text{TRalg}(\cdot)$, the parameters from Theorem 1 with input $\rho = 3$, and the following length for Phase A (resp. with ι_2, c_2)*

$$T_1 = \max \left\{ \iota_1, c_1 d^N (dr + r^N) \omega^{-2}, c_1^{\frac{3}{5}} d^{\frac{3N}{5}} (dr + r^N)^{\frac{3}{5}} \omega^{-\frac{6}{5}} T^{\frac{2}{5}} \right\},$$

with probability at least $1 - \delta - \delta_1$ (resp. $1 - \delta - \delta_2$), the regret of TOFU can be bounded as

$$R(T) \leq \tilde{O}\left(CT_1 + d^2 r^{N-2} \sqrt{T}\right)$$

The above corollary adopts $\rho = 3$, i.e., the norm constraint in Eqn. (8) is on blocks with at least three tails. This choice

is conscious with respect to the function $\eta(T_1)$ from Examples 1 and 2 as it lays aside as many parameters as possible without letting them negatively impact the bandit learning. In particular, with this choice, the length T_1 can be optimized as in Corollary 1 (which is of order $O(T^{2/5})$) and thus the dominating term (regarding the T -dependency) of the regret in Corollary 1 is the last one of order $\tilde{O}(d^2 r^{N-2} \sqrt{T})$.

This obtained regret of order $\tilde{O}(d^2 r^{N-2} \sqrt{T})$ is compared with several existing results in the following (see also Table I). First, if directly adopting linear bandits algorithms such as Lin-UCB [3] on the vectorized system, a regret of order $\tilde{O}(d^N \sqrt{T})$ would incur as the low-rank structure is not used. A second approach is to matricize the system and adopt algorithms for matrix bandits [8], [9], [21]. The state-of-the-art ESTT/ESTS [9] can then achieve a regret of order $\tilde{O}(d^{\lceil \frac{N}{2} \rceil} r^{\lfloor \frac{N}{2} \rfloor} \sqrt{T})$ (see Appendix E), which is still inefficient as matricization does not preserve all the structure information. At last, for [10] on tensor bandits, if we modify it to have general (instead of one-hot) tensors as actions, a regret of order $\tilde{O}(d^{N-1} r \sqrt{T})$ occurs as it does not fully consider the high-dimensional benefits (see Remark 1). Thus, compared with the best existing regret of order $\tilde{O}(d^{\lceil \frac{N}{2} \rceil} r^{\lfloor \frac{N}{2} \rfloor} \sqrt{T})$, TOFU has an improvement of a multiplicative factor of order $\tilde{O}((d/r)^{\lceil \frac{N}{2} \rceil - 2})$, which grows exponentially in N . Hence, this benefit becomes more significant in higher-order problems.

While TOFU improves existing results, we further compare it against the following new regret lower bound.

Theorem 2. *Assume $r^N \leq 2T$ and for all $t \in [T]$, let $\mathbb{A}_t = \mathbb{A} := \{\mathcal{A} \in \mathbb{R}^{d \times d \times \dots \times d} : \|\mathcal{A}\|_F \leq 1\}$ and ε_t be a sequence of independent standard Gaussian noise. Then, for any policy, there exists a system tensor $\mathcal{X} \in \mathbb{R}^{d \times d \times \dots \times d}$ with a multilinear rank (r, r, \dots, r) and $\|\mathcal{X}\|_F^2 = O(r^{2N}/T)$ such that $\mathbb{E}_{\mathcal{X}}[R(T)] = \Omega(r^N \sqrt{T})$, where the expectation is taken with respect to the interaction of the policy and the system.*

Compared with this lower bound, TOFU is sub-optimal only up to an additional $O((d/r)^2)$ factor (which does not scale with N). We conjecture that a slightly tighter regret lower bound of order $\Omega(dr^{N-1} \sqrt{T})$ can be established, which reduces to that of $\Omega(dr \sqrt{T})$ in matrix bandits ($N = 2$) [21].

V. CONCLUSIONS

This work studied a general tensor bandits problem, where high-dimensional tensors characterize action and system parameters. Motivated by practical applications, the system tensor is modeled to be low-rank. To tackle this high-dimensional but low-rank problem, a novel algorithm named TOFU was proposed. TOFU adopts tensor regression techniques to estimate low-dimensional subspaces associated with the system tensor. The obtained estimates are then used to transform the challenging problem of low-rank tensor bandits into an equivalent but easier one of norm-constrained linear bandits. The theoretical analysis provided a regret guarantee of TOFU, which is shown to be exponentially more efficient than existing results. A novel performance lower bound was also established, further demonstrating the superiority of TOFU.

REFERENCES

- [1] T. Lattimore and C. Szepesvári, *Bandit algorithms*. Cambridge University Press, 2020.
- [2] S. Bubeck, N. Cesa-Bianchi *et al.*, “Regret analysis of stochastic and nonstochastic multi-armed bandit problems,” *Foundations and Trends® in Machine Learning*, vol. 5, no. 1, pp. 1–122, 2012.
- [3] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári, “Improved algorithms for linear stochastic bandits,” *Advances in neural information processing systems*, vol. 24, 2011.
- [4] W. Chu, L. Li, L. Reyzin, and R. Schapire, “Contextual bandits with linear payoff functions,” in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 208–214.
- [5] E. Frolov and I. Oseledets, “Tensor methods and recommender systems,” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 7, no. 3, p. e1201, 2017.
- [6] P. Symeonidis and A. Zioupos, *Matrix and tensor factorization techniques for recommender systems*. Springer, 2016, vol. 1.
- [7] E. E. Papalexakis, C. Faloutsos, and N. D. Sidiropoulos, “Tensors for data mining and data fusion: Models, applications, and scalable algorithms,” *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 8, no. 2, pp. 1–44, 2016.
- [8] K.-S. Jun, R. Willett, S. Wright, and R. Nowak, “Bilinear bandits with low-rank structure,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 3163–3172.
- [9] Y. Kang, C.-J. Hsieh, and T. C. M. Lee, “Efficient frameworks for generalized low-rank matrix bandit problems,” in *Advances in Neural Information Processing Systems*, 2022.
- [10] J. Zhou, B. Hao, Z. Wen, J. Zhang, and W. W. Sun, “Stochastic low-rank tensor bandits for multi-dimensional online decision making,” *arXiv e-prints*, pp. arXiv–2007, 2020.
- [11] T. G. Kolda and B. W. Bader, “Tensor decompositions and applications,” *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
- [12] N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos, “Tensor decomposition for signal processing and machine learning,” *IEEE Transactions on Signal Processing*, vol. 65, no. 13, pp. 3551–3582, 2017.
- [13] H. Zhou, L. Li, and H. Zhu, “Tensor regression with applications in neuroimaging data analysis,” *Journal of the American Statistical Association*, vol. 108, no. 502, pp. 540–552, 2013.
- [14] X. Li, D. Xu, H. Zhou, and L. Li, “Tucker tensor regression and neuroimaging analysis,” *Statistics in Biosciences*, vol. 10, no. 3, pp. 520–545, 2018.
- [15] B. Yaman, S. Weingärtner, N. Kargas, N. D. Sidiropoulos, and M. Akçakaya, “Low-rank tensor models for improved multidimensional mri: Application to dynamic cardiac t_1 mapping,” *IEEE transactions on computational imaging*, vol. 6, pp. 194–207, 2019.
- [16] C. I. Kanatsoulis, X. Fu, N. D. Sidiropoulos, and M. Akçakaya, “Tensor completion from regular sub-nyquist samples,” *IEEE Transactions on Signal Processing*, vol. 68, pp. 1–16, 2019.
- [17] S. Gandy, B. Recht, and I. Yamada, “Tensor completion and low-n-rank tensor recovery via convex optimization,” *Inverse problems*, vol. 27, no. 2, p. 025010, 2011.
- [18] P. Jain and S. Oh, “Provable tensor factorization with missing data,” *Advances in Neural Information Processing Systems*, vol. 27, 2014.
- [19] A. R. Zhang, Y. Luo, G. Raskutti, and M. Yuan, “Islet: Fast and optimal low-rank tensor regression via importance sketching,” *SIAM journal on mathematics of data science*, vol. 2, no. 2, pp. 444–479, 2020.
- [20] T. Ahmed, H. Raja, and W. U. Bajwa, “Tensor regression using low-rank and sparse tucker decompositions,” *SIAM Journal on Mathematics of Data Science*, vol. 2, no. 4, pp. 944–966, 2020.
- [21] Y. Lu, A. Meisami, and A. Tewari, “Low-rank generalized linear bandit problems,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 460–468.
- [22] M. Valko, R. Munos, B. Kveton, and T. Kocák, “Spectral bandits for smooth graph functions,” in *International Conference on Machine Learning*. PMLR, 2014, pp. 46–54.
- [23] R. Han, R. Willett, and A. R. Zhang, “An optimal statistical and computational framework for generalized tensor estimation,” *The Annals of Statistics*, vol. 50, no. 1, pp. 1–29, 2022.
- [24] D. Xia, M. Yuan, and C.-H. Zhang, “Statistically optimal and computationally efficient low rank tensor completion from noisy entries,” *The Annals of Statistics*, vol. 49, no. 1, 2021.
- [25] K. Jang, K.-S. Jun, S.-Y. Yun, and W. Kang, “Improved regret bounds of bilinear bandits using action space analysis,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 4744–4754.
- [26] T. Idé, K. Murugesan, D. Bounieffouf, and N. Abe, “Targeted advertising on social networks using online variational tensor regression,” *arXiv preprint arXiv:2208.10627*, 2022.

APPENDIX A
RELATED WORKS

Linear bandits. As one of the most well-studied MAB settings, linear bandits adopt vectorized features to characterize actions and system parameters. Many provably efficient algorithms have been proposed for linear bandits and achieve (nearly) minimax optimal regret, with LinUCB being the representative design [3], [4]; see [1] for a comprehensive review.

Matrix bandits. Several works extend the vectorized (one-dimensional) features in linear bandits to two dimensions [8], [9], [21], [25], i.e., matrices as features. With similar motivations as this work, this line of research on “matrix bandits” mostly focuses on leveraging the assumed low-rank property of the system matrix, and the recent work [9] achieves nearly order-optimal performance.

Tensor bandits. The concept of tensor bandits is first proposed in [10], which generalizes the problem of linear bandits to high dimensions, i.e, use tensors to characterize actions and systems. However, in [10], the action tensors are restricted to be *one-hot*. This work instead considers general action tensors and thus covers more scenarios. Moreover, this work has a different utilization of the estimated low-dimensional subspaces compared with [10]; see details in Section III and Remark 1. Lastly, instead of the Tucker decomposition adopted in [10] and this work, a recent work [26] studies the tensor bandits problem from the Canonical polyadic decomposition (CPD) perspective, whose results are however not directly comparable to this work due to different settings.

APPENDIX B

THE PROBLEM EQUIVALENCE: DERIVATION OF EQN. (3)

With

$$\begin{aligned}\hat{\mathcal{B}} &:= \mathcal{A} \times_{n \in [N]} [\hat{U}_n, \hat{U}_{n,\perp}]^\top; \\ \hat{\mathcal{Y}} &:= \mathcal{X} \times_{n \in [N]} [\hat{U}_n, \hat{U}_{n,\perp}]^\top,\end{aligned}$$

it holds that

$$\begin{aligned}\langle \hat{\mathcal{B}}, \hat{\mathcal{Y}} \rangle &= \langle \mathcal{M}_{(1)}(\hat{\mathcal{B}}), \mathcal{M}_{(1)}(\hat{\mathcal{Y}}) \rangle \\ &\stackrel{(a)}{=} \langle [\hat{U}_1, \hat{U}_{1,\perp}]^\top \mathcal{M}_{(1)}(\mathcal{A}) J^\top, [\hat{U}_1, \hat{U}_{1,\perp}]^\top \mathcal{M}_{(1)}(\mathcal{X}) J^\top \rangle \\ &= \text{tr} \left(J \mathcal{M}_{(1)}(\mathcal{A})^\top [\hat{U}_1, \hat{U}_{1,\perp}] [\hat{U}_1, \hat{U}_{1,\perp}]^\top \mathcal{M}_{(1)}(\mathcal{X}) J^\top \right) \\ &= \text{tr} \left(J \mathcal{M}_{(1)}(\mathcal{A})^\top \mathcal{M}_{(1)}(\mathcal{X}) J^\top \right) \\ &\stackrel{(b)}{=} \langle \mathcal{A} \times_{n=2}^N [\hat{U}_n, \hat{U}_{n,\perp}]^\top, \mathcal{X} \times_{n=2}^N [\hat{U}_n, \hat{U}_{n,\perp}]^\top \rangle \\ &\stackrel{(c)}{=} \langle \mathcal{A} \times_{n=3}^N [\hat{U}_n, \hat{U}_{n,\perp}]^\top, \mathcal{X} \times_{n=3}^N [\hat{U}_n, \hat{U}_{n,\perp}]^\top \rangle \\ &= \dots \\ &\stackrel{(d)}{=} \langle \mathcal{A}, \mathcal{X} \rangle\end{aligned}$$

where

$$J := [\hat{U}_N, \hat{U}_{N,\perp}]^\top \otimes [\hat{U}_{N-1}, \hat{U}_{N-1,\perp}]^\top \otimes \dots \otimes [\hat{U}_2, \hat{U}_{2,\perp}]^\top$$

with \otimes representing the Kronecker product between matrices. Equalities (a) and (b) use the property (see

[11]) that $\mathcal{Z} = \mathcal{H} \times_{n \in [N]} V_n \Leftrightarrow \mathcal{M}_{(n)}(\mathcal{Z}) = V_n \mathcal{M}_{(n)}(\mathcal{H})(V_N \otimes \dots \otimes V_{n+1} \otimes V_{n-1} \otimes \dots \otimes V_1)^\top$. Equalities (c) and (d) recursively follow similar arguments in the previous steps. Then, following basic properties regarding tensor mode product [11], it can be shown that

$$\hat{\mathcal{Y}} := \mathcal{X} \times_{n \in [N]} [\hat{U}_n, \hat{U}_{n,\perp}]^\top = \mathcal{G} \times_{n \in [N]} \left([\hat{U}_n, \hat{U}_{n,\perp}]^\top U_n \right),$$

which completes the derivation.

APPENDIX C

UPPER BOUND ANALYSIS: PROOF OF THEOREM 1

Lemma 1. *Under Assumption 2, it holds that*

$$\|\hat{U}_{n,\perp}^\top U_n\|_F \leq \frac{\eta(T_1)}{\omega_n}, \quad \forall n \in [N],$$

where $\omega_n := \omega_{\min}(\mathcal{M}_n(\mathcal{X}))$.

Proof. It holds that

$$\begin{aligned}\|\hat{\mathcal{X}} - \mathcal{X}\|_F &= \|\mathcal{M}_n(\hat{\mathcal{X}}) - \mathcal{M}_n(\mathcal{X})\|_F \\ &\stackrel{(a)}{=} \|\mathcal{M}_n(\hat{\mathcal{X}}) - U_n U_n^\top \mathcal{M}_n(\mathcal{X})\|_F \\ &\stackrel{(b)}{\geq} \|\hat{U}_{n,\perp}^\top \mathcal{M}_n(\hat{\mathcal{X}}) - \hat{U}_{n,\perp}^\top U_n U_n^\top \mathcal{M}_n(\mathcal{X})\|_F \\ &\stackrel{(c)}{=} \|\hat{U}_{n,\perp}^\top U_n U_n^\top \mathcal{M}_n(\mathcal{X})\|_F \\ &\stackrel{(d)}{\geq} \omega_{\min}(U_n^\top \mathcal{M}_n(\mathcal{X})) \|\hat{U}_{n,\perp}^\top U_n\|_F \\ &= \omega_n \|\hat{U}_{n,\perp}^\top U_n\|_F,\end{aligned}$$

where equality (a) is because U_n contains the left singular vectors of $\mathcal{M}_n(\mathcal{X})$; inequality (b) is from the fact that for a matrix U with orthonormal columns and an arbitrary compatible matrix X , it holds that $\|X\|_F \geq \|U^\top X\|_F$; inequality (c) uses the observation that $\hat{U}_{n,\perp}^\top \hat{U}_n$ is a null matrix; inequality (d) is from the fact that for any compatible matrices X and Y , $\|XY\|_F \geq \min\{\omega_{\min}(X)\|Y\|_F, \omega_{\min}(Y)\|X\|_F\}$. The lemma is proved by plugging Assumption 2 into the above inequality. \square

Lemma 2. *Under Assumption 2, the norm of the following block with k tails in $\hat{\mathcal{Y}}$ can be bounded as:*

$$\left\| \underbrace{\hat{\mathcal{Y}}_{:r; :r, \dots, :r, r+1; :r+1, \dots, :r+1}}_{N-k \text{ modes}} \right\|_F \leq \frac{C(\eta(T_1))^k}{\omega^k},$$

and this bound symmetrically holds for all $\binom{N}{k}$ blocks with k tails.

Proof. It holds that

$$\begin{aligned}&\left\| \hat{\mathcal{Y}}_{:r; :r, \dots, :r, r+1; :r+1, \dots, :r+1} \right\|_F \\ &= \left\| \mathcal{G} \times_{n \in [N-k]} (\hat{U}_n^\top U_n) \times_{n' \in [N-k+1:N]} (\hat{U}_{n',\perp}^\top U_{n'}) \right\|_F \\ &\stackrel{(a)}{\leq} \|\mathcal{G}\|_F \prod_{n \in [N-k]} \left\| (\hat{U}_n^\top U_n) \right\|_2 \prod_{n' \in [N-k+1:N]} \left\| (\hat{U}_{n',\perp}^\top U_{n'}) \right\|_2 \\ &\stackrel{(b)}{\leq} \frac{\|\mathcal{G}\|_F \cdot (\eta(T_1))^k}{\omega^k} \leq \frac{C(\eta(T_1))^k}{\omega^k},\end{aligned}$$

APPENDIX D

UPPER BOUND ANALYSIS: PROOF OF COROLLARY 1

where inequality (a) repeatedly uses the fact that for any arbitrary matrices X and Y , it holds that $\|XY\|_F \leq \min\{\|X\|_2\|Y\|_F, \|X\|_F\|Y\|_2\}$; inequality (b) utilizes Lemma 1, Assumption 3, and the fact that for two compatible matrices U and V with orthonormal columns, it holds that $\|V^\top U\|_2 \leq 1$. \square

Lemma 3 (Corollary 1 of [8]). *If Eqn. (8) holds, with*

$$\lambda_\perp = \frac{T}{q(\rho) \log(1 + T/\lambda)}$$

the regret of LowOFUL (adopted in Phase B) for T steps is, with probability at least $1 - \delta$, bounded by

$$\tilde{O}\left(\left(q(\rho) + \sqrt{q(\rho)\lambda}C + \sqrt{TC}_\perp\right)\sqrt{T}\right).$$

Proof. Detailed proofs can be found in [8]. \square

Then, the proof of Theorem 1 is presented in the following.

Proof of Theorem 1. For Phase A with length T_1 , its regret can be bounded as

$$R_A(T) \leq 2CT_1,$$

since the mean rewards are bounded between $[-C, C]$ with $\|\mathcal{A}_t\|_F \leq 1$ and $\|\mathcal{X}\|_F \leq C$.

After Phase A, based on Lemma 2, we have that for all $k \in [N]$, it holds that

$$\|\hat{\mathbf{y}}_{q(k)+1:q(k+1)}\|_F^2 \leq \binom{N}{k} \frac{C^2(\eta(T_1))^{2k}}{\omega^{2k}},$$

Thus, with $q(\rho)$ in Eqn. (6), it can be further shown that

$$\begin{aligned} \|\hat{\mathbf{y}}_{q(\rho)+1:d^N}\|_F^2 &= \sum_{k \in [\rho:N]} \|\hat{\mathbf{y}}_{q(k)+1:q(k+1)}\|_F^2 \\ &\leq \sum_{k \in [\rho:N]} \binom{N}{k} \frac{C^2(\eta(T_1))^{2k}}{\omega^{2k}} \\ &\leq \frac{2^N \cdot C^2(\eta(T_1))^{2\rho}}{\omega^{2\rho}}, \end{aligned}$$

where the last inequality uses the condition that $\eta(T_1) \leq \omega$. This inequality validates Eqn. (8) with the parameter $C_\perp = 2^{N/2}C(\eta(T_1))^\rho\omega^{-\rho}$ in Theorem 1.

Then, with the specified parameter, according to Lemma 3, Phase B would lead to a regret bounded as

$$\begin{aligned} R_B(T) &\leq \tilde{O}\left(\left(q(\rho) + \sqrt{q(\rho)\lambda}C + \sqrt{TC}_\perp\right)\sqrt{T}\right) \\ &= \tilde{O}\left(d^{\rho-1}r^{N-\rho+1}\sqrt{T} + \frac{C(\eta(T_1))^\rho}{\omega^\rho}T\right), \end{aligned}$$

where the last step uses the facts that $q(\rho) = O(2^N d^{\rho-1} r^{N-\rho+1})$ and $N = O(1)$ (thus $2^N = O(1)$). Thus, the overall regret guarantee can be obtained as

$$\begin{aligned} R(T) &= R_A(T) + R_B(T) \\ &\leq \tilde{O}\left(CT_1 + d^{\rho-1}r^{N-\rho+1}\sqrt{T} + \frac{C(\eta(T_1))^\rho}{\omega^\rho}T\right), \end{aligned}$$

which concludes the proof. \square

Proof. In the following, we prove the case for Example 1. The proof for Example 2 can be similarly constructed. First, with probability at least $1 - \delta - \delta_1$, it simultaneously holds that

$$R(T) \leq \tilde{O}\left(CT_1 + d^{\rho-1}r^{N-\rho+1}\sqrt{T} + \frac{C(\eta(T_1))^\rho}{\omega^\rho}T\right)$$

and

$$\|\hat{\mathcal{X}} - \mathcal{X}\|_F^2 \leq \eta(T_1) = \sqrt{\frac{c_1 d^N (dr + r^N)}{T_1}}.$$

Thus, if $\rho = 3$ as specified in Corollary 1, it holds that

$$R(T) \leq \tilde{O}\left(CT_1 + d^2 r^{N-2} \sqrt{T} + \frac{Cc_1^{\frac{3}{2}} d^{\frac{3N}{2}} (dr + r^N)^{\frac{3}{2}}}{\omega^3 T_1^{\frac{3}{2}}} T\right).$$

With the following choice of

$$T_1 = \max\left\{\iota_1, \frac{c_1 d^N (dr + r^N)}{\omega^2}, \frac{c_1^{\frac{3}{2}} d^{\frac{3N}{2}} (dr + r^N)^{\frac{3}{2}}}{\omega^{\frac{6}{5}}} T^{\frac{2}{5}}\right\},$$

the threshold requirement in Example 1 can be satisfied (i.e., $T_1 \geq \iota_1$) and it can be verified that $\eta(T_1) \leq \omega$. Thus, with probability at least $1 - \delta - \delta_1$, the regret can be bounded as

$$\begin{aligned} R(T) &\leq \tilde{O}\left(CT_1 + d^2 r^{N-2} \sqrt{T} + \frac{Cc_1^{\frac{3}{2}} d^{\frac{3N}{2}} (dr + r^N)^{\frac{3}{2}}}{\omega^3 T_1^{\frac{3}{2}}} T\right) \\ &= \tilde{O}\left(CT_1 + d^2 r^{N-2} \sqrt{T}\right), \end{aligned}$$

where the selected value of T_1 is adopted. The proof is then concluded. \square

APPENDIX E

REGRET OF MATRICIZED ESTT/ESTS

ESTT/ESTS [9] deals with the matrix bandits problem where the actions and system parameters are characterized by matrices. In particular, when the system matrix is of size (D_1, D_2) and rank R , a regret of $\tilde{O}((D_1 + D_2)R\sqrt{T})$ can be obtained. The straightforward way to matricize the order- N system tensor considered in this work is along one mode, e.g., as $\mathcal{M}_n(X)$. This obtained system matrix $\mathcal{M}_n(X)$ would be of size (d, d^{N-1}) and rank r , which results in a regret of $\tilde{O}(d^{N-1}r\sqrt{T})$ with ESTT/ESTS. However, if we combine $\lfloor N/2 \rfloor$ modes in the system tensor in one matrix dimension (e.g., row), and the remaining $\lfloor N/2 \rfloor$ modes in the other matrix dimension (e.g., column), a matrix of size $(d^{\lfloor N/2 \rfloor}, d^{\lfloor N/2 \rfloor})$ can be obtained with rank $r^{\lfloor N/2 \rfloor}$. Using this matricization, ESTT/ESTS can obtain a regret of $\tilde{O}(d^{\lfloor N/2 \rfloor} r^{\lfloor N/2 \rfloor} \sqrt{T})$, which is much better than $\tilde{O}(d^{N-1}r\sqrt{T})$ and thus adopted as the baseline result in the main paper.

APPENDIX F

LOWER BOUND ANALYSIS: PROOF OF THEOREM 2

Proof of Theorem 2. In the following, for $i = (i_1, i_2, \dots, i_N)$, we adopt the simplified notation that

$$\mathcal{X}_i := \mathcal{X}_{i_1, i_2, \dots, i_N}.$$

With $\Delta := \frac{1}{8\sqrt{3}}\sqrt{\frac{r^N}{T}}$, we design $\mathfrak{G} = \{\mathcal{G} \text{ that satisfies Eqn. (9)}\} \subseteq \mathbb{R}^{r \times r \times \dots \times r}$ and $\mathfrak{X} = \{\mathcal{X} = \mathcal{G} \times_{n \in [N]} U_n : \mathcal{G} \in \mathfrak{G}\} \subseteq \mathbb{R}^{d \times d \times \dots \times d}$, where

$$\mathcal{G}_i \in \{\pm\Delta\}, \quad \forall i = (i_1, \dots, i_N) \in [d] \times \dots \times [d] \quad (9)$$

and

$$U_n = \begin{bmatrix} I_r \\ 0_{(d-r) \times r} \end{bmatrix} \in \mathbb{R}^{d \times r}, \quad \forall n \in [N].$$

It can be noted that each $\mathcal{X} \in \mathfrak{X}$ has a multi-linear rank at most (r, r, \dots, r) . For $i = (i_1, i_2, \dots, i_N) \in [r] \times [r] \times \dots \times [r]$, we define

$$\tau_i = T \wedge \min \left\{ t : \sum_{\tau \in [t]} \mathcal{A}_{t;i}^2 \geq \frac{T}{r^N} \right\}.$$

For a fixed $\mathcal{X} \in \mathfrak{X}$, we have

$$\begin{aligned} \mathbb{E}_{\mathcal{X}}[R(T)] &= \mathbb{E}_{\mathcal{X}} \left[\sum_{t \in [T]} \langle \mathcal{A}^* - \mathcal{A}_t, \mathcal{X} \rangle \right] \\ &= \Delta \mathbb{E}_{\mathcal{X}} \left[\sum_{t \in [T]} \sum_{i \in [r] \times \dots \times [r]} \left(\frac{1}{r^{\frac{N}{2}}} - \mathcal{A}_{t;i} \cdot \text{sign}(\mathcal{X}_i) \right) \right] \\ &\geq \frac{\Delta r^{\frac{N}{2}}}{2} \mathbb{E}_{\mathcal{X}} \left[\sum_{t \in [T]} \sum_{i \in [r] \times \dots \times [r]} \left(\frac{1}{r^{\frac{N}{2}}} - \mathcal{A}_{t;i} \cdot \text{sign}(\mathcal{X}_i) \right)^2 \right] \\ &\geq \frac{\Delta r^{\frac{N}{2}}}{2} \sum_{i \in [r] \times \dots \times [r]} \mathbb{E}_{\mathcal{X}} \left[\sum_{t \in [\tau_i]} \left(\frac{1}{r^{\frac{N}{2}}} - \mathcal{A}_{t;i} \cdot \text{sign}(\mathcal{X}_i) \right)^2 \right], \end{aligned}$$

where the first inequality uses the fact that $\|\mathcal{A}_t\|_F \leq 1$.

Let $\mathcal{G}' \in \mathfrak{G}$ be another tensor such that $\mathcal{G}' = \mathcal{G}$ except $\mathcal{G}'_i = -\mathcal{G}_i$ with $i = (i_1, i_2, \dots, i_N)$. Then, with $\mathcal{X}' = \mathcal{G}' \times_{n \in [N]} U_n$, it also holds that $\mathcal{X}' = \mathcal{X}$ except $\mathcal{X}'_i = -\mathcal{X}_i$. For $x \in \{\pm 1\}$, we define

$$\kappa_i(x) = \sum_{t \in [\tau_i]} \left(\frac{1}{r^{\frac{N}{2}}} - \mathcal{A}_{t;i} \cdot x \right)^2.$$

Let \mathbb{P} and \mathbb{P}' be the distributions of $\kappa_i(1)$ with respect to the player interaction measure induced by \mathcal{X} and \mathcal{X}' , respectively. Then, it holds that

$$\begin{aligned} \mathbb{E}_{\mathcal{X}}[\kappa_i(1)] &\stackrel{(a)}{\geq} \mathbb{E}_{\mathcal{X}'}[\kappa_i(1)] - \left(\frac{4T}{r^N} + 2 \right) \sqrt{\frac{1}{2} D(\mathbb{P}, \mathbb{P}')} \\ &\stackrel{(b)}{\geq} \mathbb{E}_{\mathcal{X}'}[\kappa_i(1)] - \left(\frac{4T}{r^N} + 2 \right) \Delta \sqrt{\mathbb{E}_{\mathcal{X}} \left[\sum_{t \in [\tau_i]} \mathcal{A}_{t;i}^2 \right]} \end{aligned}$$

$$\begin{aligned} &\geq \mathbb{E}_{\mathcal{X}'}[\kappa_i(1)] - \left(\frac{4T}{r^N} + 2 \right) \Delta \sqrt{\frac{T}{r^N} + 1} \\ &\stackrel{(c)}{\geq} \mathbb{E}_{\mathcal{X}'}[\kappa_i(1)] - \frac{8\sqrt{3}T\Delta}{r^N} \sqrt{\frac{T}{r^N}} \end{aligned}$$

where $D(\cdot, \cdot)$ is the relative entropy between two probability measures. Inequality (a) uses the result in Exercise 14.4 of [1], the Pinsker's inequality, and the bound

$$\kappa_i(1) \leq 2 \sum_{t \in [\tau_i]} \frac{1}{r^N} + 2 \sum_{t \in [\tau_i]} \mathcal{A}_{t;i}^2 \leq \frac{4T}{r^N} + 2,$$

where the definition of τ_i is used for the last inequality. Inequality (b) is from the chain rule for the relative entropy up to a stopping time in Exercise 15.7 of [1] as follows:

$$\begin{aligned} D(\mathbb{P}, \mathbb{P}') &\leq \frac{1}{2} \mathbb{E}_{\mathcal{X}} \left[\sum_{t \in [\tau_i]} \langle \mathcal{A}_t, \mathcal{X} - \mathcal{X}' \rangle^2 \right] \\ &= 2\Delta^2 \mathbb{E}_{\mathcal{X}} \left[\sum_{t \in [\tau_i]} \mathcal{A}_{t;i}^2 \right]. \end{aligned}$$

Inequality (c) is from the assumption that $r^N \leq 2T$.

Then, it holds that

$$\begin{aligned} &\mathbb{E}_{\mathcal{X}}[\kappa_i(1)] + \mathbb{E}_{\mathcal{X}'}[\kappa_i(-1)] \\ &\geq \mathbb{E}_{\mathcal{X}'}[\kappa_i(1) + \kappa_i(-1)] - \frac{8\sqrt{3}T\Delta}{r^N} \sqrt{\frac{T}{r^N}} \\ &= 2\mathbb{E}_{\mathcal{X}'} \left[\frac{\tau_i}{r^N} + \sum_{t \in [\tau_i]} \mathcal{A}_{t;i}^2 \right] - \frac{8\sqrt{3}T\Delta}{r^N} \sqrt{\frac{T}{r^N}} \\ &\geq \frac{2T}{r^N} - \frac{8\sqrt{3}T\Delta}{r^N} \sqrt{\frac{T}{r^N}} \\ &\geq \frac{T}{r^N}, \end{aligned}$$

where the last inequality uses the definition of Δ .

The proof is completed using an averaging number argument on the following quantity:

$$\begin{aligned} \sum_{\mathcal{X} \in \mathfrak{X}} \mathbb{E}_{\mathcal{X}}[R(T)] &\geq \frac{\Delta r^{\frac{N}{2}}}{2} \sum_{i \in [r] \times \dots \times [r]} \sum_{\mathcal{X} \in \mathfrak{X}} \mathbb{E}_{\mathcal{X}}[\kappa_i(\text{sign}(\mathcal{X}_i))] \\ &= \frac{\Delta r^{\frac{N}{2}}}{2} \sum_{i \in [r] \times \dots \times [r]} \sum_{\mathcal{X}/\mathcal{X}_i \in \{\pm\Delta\}^{r^N-1}} \sum_{\mathcal{X}_i \in \{\pm\Delta\}} \mathbb{E}_{\mathcal{X}}[\kappa_i(\text{sign}(\mathcal{X}_i))] \\ &\geq \frac{\Delta r^{\frac{N}{2}}}{2} \sum_{i \in [r] \times \dots \times [r]} \sum_{\mathcal{X}/\mathcal{X}_i \in \{\pm\Delta\}^{r^N-1}} \frac{T}{r^N} \\ &= 2^{r^N-2} \Delta r^{\frac{N}{2}} T. \end{aligned}$$

Hence there exists $\mathcal{G} \in \mathfrak{G}$ and $\mathcal{X} = \mathcal{G} \times_{n \in [N]} U_n$ such that

$$\mathbb{E}_{\mathcal{X}}[R(T)] \geq \frac{\Delta r^{\frac{N}{2}} T}{4} = \frac{r^N \sqrt{T}}{32\sqrt{3}},$$

which concludes the proof. \square

APPENDIX G
DISCUSSION AND FUTURE DIRECTIONS

Tighter upper and lower regret bounds. As mentioned at the end of Section IV, it is conjectured that a slightly tighter regret lower bound of order $\Omega(dr^{N-1}\sqrt{T})$ exists, which reduces to $\Omega(dr\sqrt{T})$ for matrix bandits ($N = 2$) [21]. It would be an interesting question to (dis)prove this conjecture, and we hope the proof of the current Theorem 2 can be inspiring. On the other hand, it remains a challenging problem to further tighten the performance upper bound established in Theorem 1 and Corollary 1. Inspired by the recent success in matrix bandits [9], one potential direction is to only estimate the subspaces $(\hat{U}_1, \dots, \hat{U}_N)$ in Phase A, because $\hat{\mathcal{G}}$ is not used in later learning. In particular, it would be sufficient to obtain an estimate $\hat{\mathcal{X}}$ that approaches $v\mathcal{X}$ with v as an unknown constant. If corresponding techniques can be proposed in the study of tensor estimation, the general framework of TOFU can be smoothly adapted and sharper upper bounds can be similarly obtained.

Structure action sets. This work mainly investigates the problem with a low-rank tensor for system parameters. It would be valuable to also consider structured tensors for actions. Preliminary results for matrix bandits can be found in [25], where low-rank action matrices are studied.

From Tucker to CPD. Besides the Tucker decomposition, another well-known tensor decomposition is CPD. It would be interesting to study the problem of tensor bandits with a low-rank CPD, which might be able to eliminate the exponential dependency on N . However, this direction is challenging as the projections constructed in TOFU cannot be performed under the CPD formulation, and requires further investigations.