

# Contextual Online Learning for Multimedia Content Aggregation

Cem Tekin, *Member, IEEE*, Mihaela van der Schaar, *Fellow, IEEE*

**Abstract**—The last decade has witnessed a tremendous growth in the volume as well as the diversity of multimedia content generated by a multitude of sources (news agencies, social media, etc.). Faced with a variety of content choices, consumers are exhibiting diverse preferences for content; their preferences often depend on the context in which they consume content as well as various exogenous events. To satisfy the consumers’ demand for such diverse content, multimedia content aggregators (CAs) have emerged which gather content from numerous multimedia sources. A key challenge for such systems is to accurately predict what type of content each of its consumers prefers in a certain context, and adapt these predictions to the evolving consumers’ preferences, contexts and content characteristics. We propose a novel, distributed, online multimedia content aggregation framework, which gathers content generated by multiple heterogeneous producers to fulfill its consumers’ demand for content. Since both the multimedia content characteristics and the consumers’ preferences and contexts are unknown, the optimal content aggregation strategy is unknown a priori. Our proposed content aggregation algorithm is able to learn online what content to gather and how to match content and users by exploiting similarities between consumer types. We prove bounds for our proposed learning algorithms that guarantee both the accuracy of the predictions as well as the learning speed. Importantly, our algorithms operate efficiently even when feedback from consumers is missing or content and preferences evolve over time. Illustrative results highlight the merits of the proposed content aggregation system in a variety of settings.

**Index Terms**—Social multimedia, distributed online learning, content aggregation, multi-armed bandits.

## I. INTRODUCTION

A plethora of multimedia applications (web-based TV [2], [3], personalized video retrieval [4], personalized news aggregation [5], etc.) are emerging which require matching multimedia content generated by distributed sources with consumers exhibiting different interests. The matching is often performed by CAs (e.g., Dailymotion, Metacafe [6]) that are responsible for mining the content of numerous multimedia sources in search of finding content which is interesting for the users. Both the characteristics of the content and preference of the consumers are evolving over time. An example of the system with users, CAs and multimedia sources is given in Fig. 1.

Each user is characterized by its context, which is a real-valued vector, that provides information about the users’

Copyright (c) 2015 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

This work is partially supported by the grants NSF CNS 1016081 and AFOSR DDDAS.

C. Tekin and Mihaela van der Schaar are in Department of Electrical Engineering, UCLA, Los Angeles, CA, 90095. Email: cmtkn@ucla.edu, mihaela@ee.ucla.edu.

This online technical report is an extended version of the paper that appeared in IEEE Transactions on Multimedia [1].

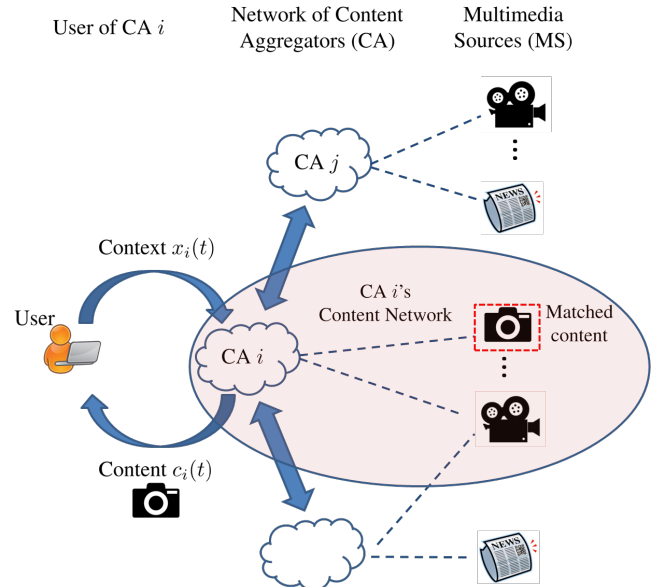


Fig. 1. Operation of the distributed content aggregation system. (i) A user with type/context  $x_i(t)$  arrives to Content Aggregator (CA)  $i$ , (ii) CA  $i$  chooses a *matching action* (either requests content from another CA or requests content from a multimedia source in its own network).

content preferences. We assume a model where users arrive sequentially to a CA, and based on the type (context) of the user, the CA requests content from either one of the multimedia sources that it is connected to or from another CA that it is connected to. The context can represent information such as age, gender, search query, previously consumed content, etc. It may also represent the type of the device that the user is using [7] (e.g., PDA, PC, mobile phone). The CA’s role is to match its user with the most suitable content, which can be accomplished by requesting content from the most suitable multimedia source.<sup>1</sup> Since both the content generated by the multimedia sources and the user’s characteristics change over time, it is unknown to the CA which multimedia source to match with the user. This problem can be formulated as an online learning problem, where the CA learns the best matching by exploring matchings of users with different content providers. After a particular content matching is made, the user “consumes” the content, and provides feedback/rating, such as *like* or *dislike*.<sup>2</sup> It is this feedback that helps a CA learn the preferences of its users and the characteristics of the content that is provided by the multimedia sources. Since this

<sup>1</sup>Although we use the term *request* to explain how content from a multimedia source is mined, our proposed method works also when a CA extracts the content from the multimedia source, without any decision making performed by the multimedia source.

<sup>2</sup>Our framework also works when the feedback is missing for some users.

is a learning problem we equivalently call a CA, a content learner or simply, a *learner*.

Two possible real-world applications of content aggregation are business news aggregation and music aggregation. Business news aggregators can collect information from a variety of multinational and multilingual sources and make personalized recommendations to specific individuals/companies based on their unique needs (see e.g. [8]). Music aggregators enable matching listeners with music content they enjoy both within the content network of the listeners as well as outside this network. For instance, distributed music aggregators can facilitate the sharing of music collections owned by diverse users without the need for centralized content manager/moderator/providers (see e.g. [9]). A discussion of how these applications can be modeled using our framework is given in Section III. Moreover, our proposed methods are tested on real-world datasets related to news aggregation and music aggregation in Section VII.

For each CA  $i$ , there are two types of users: *direct* and *indirect*. Direct users are the users that visit the website of CA  $i$  to search for content. *Indirect* users are the users of another CA that requests content from CA  $i$ . A CA's goal is to maximize the number of likes received from its users (both direct and indirect). This objective can be achieved by all CAs by the following *distributed* learning method: all CAs learn online which *matching action* to take for its current user, i.e., obtain content from a multimedia source that is directly connected, or request content from another CA. However, it is not trivial how to use the past information collected by the CAs in an efficient way, due to the vast number of contexts (different user types) and dynamically changing user and content characteristics. For instance, a certain type of content may become popular among users at a certain point in time, which will require the CA to obtain content from the multimedia source that generates that type of content.

To jointly optimize the performance of the multimedia content aggregation system, we propose an online learning methodology that builds on contextual bandits [10], [11]. The performance of the proposed methodology is evaluated using the notion of regret: the difference between the expected total reward (number of content likes minus costs of obtaining the content) of the best content matching strategy given complete knowledge about the user preferences and content characteristics and the expected total reward of the algorithm used by the CAs. When the user preferences and content characteristics are static, our proposed algorithms achieve sublinear regret in the number of users that have arrived to the system.<sup>3</sup> When the user preferences and content characteristics are slowly changing over time, our proposed algorithms achieve  $\epsilon$  time-averaged regret, where  $\epsilon > 0$  depends on the rate of change of the user and content characteristics.

The remainder of the paper is organized as follows. In Section II, we describe the related work and highlight the differences from our work. In Section III, we describe the decentralized content aggregation problem, the optimal content

matching scheme given the complete system model, and the regret of a learning algorithm with respect to the optimal content matching scheme. Then, we consider the model with unknown, static user preferences and content characteristics and propose a distributed online learning algorithm in Section IV. The analysis of the unknown, dynamic user preferences and content characteristics are given in Section VI. Using real-world datasets, we provide numerical results on the performance of our distributed online learning algorithms in Section VII. Finally, the concluding remarks are given in Section VIII.

## II. RELATED WORK

Related work can be categorized into two: related work on recommender systems and related work on online learning methods called *multi-armed bandits*.

### A. Related work on recommender systems and content matching

A recommender system recommends items to its users based on the characteristics of the users and the items. The goal of a recommender system is to learn which users like which items, and recommend items such that the number of likes is maximized. For instance, in [5], [12] a recommender system that learns the preferences of its users in an online way based on the ratings submitted by the users is provided. It is assumed that the true relevance score of an item for a user is a linear function of the context of the user and the features of the item. Under this assumption, an online learning algorithm is proposed. In contrast, we consider a different model, where the relevance score need not be linear in the context. Moreover, due to the distributed nature of the problem we consider, our online learning algorithms need an additional phase called the training phase, which accounts for the fact that the CAs are uncertain about the information of the other aggregators that they are linked with. We focus on the long run performance and show that the regret per unit time approaches zero when the user and content characteristics are static. An online learning algorithm for a centralized recommender which updates its recommendations as both the preferences of the users and the characteristics of items change over time is proposed in [13].

The general framework which exploits the similarities between the past users and the current user to recommend content to the current user is called collaborative filtering [14]–[16]. These methods find the similarities between the current user and the past users by examining their search and feedback patterns, and then based on the interactions with the past *similar* users, matches the user with the content that has the highest estimated relevance score. For example, the most relevant content can be the content that is liked the highest number of times by similar users. Groups of similar users can be created by various methods such as clustering [15], and then, the matching will be made based on the content matched with the past users that are in the same group.

The most striking difference between our content matching system and previously proposed is that in prior works, there is a central CA which knows the entire set of different types of

<sup>3</sup>We use index  $t$  to denote the number of users that have arrived so far. We also call  $t$  the time index, and assume that one user arrives at each time step.

	Our work	[5], [12]	[15]	[14]	[19]
Distributed	Yes	No	No	No	No
Reward model	Hölder	Linear	N/A	N/A	N/A
Confidence bounds	Yes	No	No	No	No
Regret bound	Yes	Yes	No	No	No
Dynamic user /content distribution	Yes	No	Yes	Yes	Yes

TABLE I  
COMPARISON OF OUR WORK WITH OTHER WORK IN RECOMMENDER SYSTEMS

content, and all the users arrive to this central CA. In contrast, we consider a decentralized system consisting of many CAs, many multimedia sources that these CAs are connected to, and heterogeneous user arrivals to these CAs. These CAs are cooperating with each other by only knowing the connections with their own neighbors but not the entire network topology. Hence, a CA does not know which multimedia sources another CA is connected to, but it learns over time whether that CA has access to content that the users like or not. Thus, our model can be viewed as a giant collection of individual CAs that are running in parallel.

Another line of work [17], [18] uses social streams mined in one domain, e.g., Twitter, to build a topic space that relates these streams to content in the multimedia domain. For example, in [17], Tweet streams are used to provide video recommendations in a commercial video search engine. A content adaptation method is proposed in [7] which enables the users with different types of contexts and devices to receive content that is in a suitable format to be accessed. Video popularity prediction is studied in [18], where the goal is to predict if a video will become popular in the multimedia domain, by detecting social trends in another social media domain (such as Twitter), and transferring this knowledge to the multimedia domain. Although these methods are very different from our methods, the idea of transferring knowledge from one multimedia domain to another can be carried out by CAs specialized in specific types of cross-domain content matching. For instance, one CA may transfer knowledge from tweets to predict the content which will have a high relevance/popularity for a user with a particular context, while another CA may scan through the Facebook posts of the user’s friends to calculate the context of the domain in addition to the context of the user, and provide a matching according to this.

The advantages of our proposed approach over prior work in recommender systems are: (i) systematic analysis of recommendations’ performance, including confidence bounds on the accuracy of the recommendations; (ii) no need for a priori knowledge of the users’ preferences (i.e., system learns on-the-fly); (iii) achieve high accuracy even when the users’ characteristics and content characteristics are changing over time; (iv) all these features are enabled in a network of distributed CAs.

The differences of our work from the prior work in recommender systems is summarized in Table I.

### B. Related Work on Multi-armed Bandits

Other than distributed content recommendation, our learning framework can be applied to any problem that can be formulated as a decentralized contextual bandit problem. Contextual

bandits have been studied before in [10], [11], [20]–[22] in a single agent setting, where the agent sequentially chooses from a set of alternatives with unknown rewards, and the rewards depend on the context information provided to the agent at each time step. In [5], a contextual bandit algorithm named LinUCB is proposed for recommending personalized news articles, which is variant of the UCB algorithm [23] designed for linear payoffs. Numerical results on real-world Internet data are provided, but no theoretical results on the resulting regret are derived. The main difference of our work from single agent contextual bandits is that: (i) a three phase learning algorithm with *training*, *exploration* and *exploitation* phases is needed instead of the standard two phase, i.e., *exploration* and *exploitation* phases, algorithms used in centralized contextual bandit problems; (ii) the adaptive partitions of the context space should be formed in a way that each learner/aggregator can efficiently utilize what is learned by other learners about the same context; (iii) the algorithm is robust to missing feedback (some users do not rate the content).

### III. PROBLEM FORMULATION

The system model is shown in Fig. 1. There are  $M$  content aggregators (CAs) which are indexed by the set  $\mathcal{M} := \{1, 2, \dots, M\}$ . We also call each CA a *learner* since it needs to learn which type of content to provide to its users. Let  $\mathcal{M}_{-i} := \mathcal{M} - \{i\}$  be the set of CAs that CA  $i$  can choose from to request content. Each CA has access to the contents over its content network as shown in Fig. 1. The set of contents in CA  $i$ ’s content network is denoted by  $\mathcal{C}_i$ . The set of all contents is denoted by  $\mathcal{C} := \cup_{i \in \mathcal{M}} \mathcal{C}_i$ . The system works in a discrete time setting  $t = 1, 2, \dots, T$ , where the following events happen sequentially, in each time slot: (i) a user with context  $x_i(t)$  arrives to each CA  $i \in \mathcal{M}$ ,<sup>4</sup> (ii) based on the context of its user each CA matches its user with a content (either from its own content network or by requesting content from another CA), (iii) the user provides a feedback, denoted by  $y_i(t)$ , which is either *like* ( $y_i(t) = 1$ ) or *dislike* ( $y_i(t) = 0$ ).

The set of *content matching actions* of CA  $i$  is denoted by  $\mathcal{K}_i := \mathcal{C}_i \cup \mathcal{M}_{-i}$ . Let  $\mathcal{X} = [0, 1]^d$  be the context space,<sup>5</sup> where  $d$  is the dimension of the context space. The context can include many properties of the user such as age, gender, income, previously liked content, etc. We assume that all these quantities are mapped into  $[0, 1]^d$ . For instance, this mapping can be established by feature extraction methods such as the one given in [5]. Another method is to represent each property of a user by a real number between  $[0, 1]$  (e.g., normalize the age by a maximum possible age, represent gender by set  $\{0, 1\}$ , etc.), without feature extraction. The feedback set of a user is denoted by  $\mathcal{Y} := \{0, 1\}$ . Let  $C_{\max} := \max_{i \in \mathcal{M}} |\mathcal{C}_i|$ . We assume that all CAs know  $C_{\max}$  but they do not need to know the content networks of other CAs.

The following two examples demonstrate how business news aggregation and music aggregation fits our problem formulation.

<sup>4</sup>Although in this model user arrivals are synchronous, our framework will work for asynchronous user arrivals as well.

<sup>5</sup>In general, our results will hold for any bounded subspace of  $\mathbb{R}^n$ .

**Example 1: Business news aggregation.** Consider a distributed set of news aggregators that operate in different countries (for instance a European news aggregator network as in [8]). Each news aggregator’s content network (as portrayed in Fig. 1 of the manuscript) consists of content producers (multimedia sources) that are located in specific regions/countries. Consider a user with context  $x$  (e.g. age, gender, nationality, profession) who subscribes to the CA  $A$ , which is located in the country where the user lives. This CA has access to content from local producers in that country but it can also request content from other CAs, located in different countries. Hence, a CA has access to (local) content generated in other countries. In such scenarios, our proposed system is able to recommend to the user subscribing to CA  $A$  also content from other CAs, by discovering the content that is most relevant to that user (based on its context  $x$ ) across the entire network of CAs. For instance, for a user doing business in the transportation industry, our content aggregator system may learn to recommend road construction news, accidents or gas prices from particular regions that are on the route of the transportation network of the user.

**Example 2: Music aggregation.** Consider a distributed set of music aggregators that are specialized in specific genres of music: classical, jazz, rock, rap, etc. Our proposed model allows music aggregators to share content to provide personalized recommendation for a specific user. For instance, a user that subscribes (frequents/listens) to the classical music aggregator may also like specific jazz tracks. Our proposed system is able to discover and recommend to that user also other music that it will enjoy in addition to the music available to/owned by in aggregator to which it subscribes.

#### A. User and Content Characteristics

In this paper we consider two types of user and content characteristics. First, we consider the case when the user and content characteristics are *static*, i.e., they do not change over time. For this case, for a user with context  $x$ ,  $\pi_c(x)$  denotes the probability that the user will like content  $c$ . We call this the *relevance score* of content  $c$ .

The second case we consider corresponds to the scenario when the characteristics of the users and content are *dynamic*. For online multimedia content, especially for social media, it is known that both the user and the content characteristics are dynamic and noisy [24], hence the problem exhibits concept drift [25]. Formally, a concept is the distribution of the problem, i.e., the joint distribution of the user and content characteristics, at a certain point of time [26]. Concept drift is a change in this distribution. For the case with concept drift, we propose a learning algorithm that takes into account the speed of the drift to decide what window of past observations to use in estimating the relevance score. The proposed learning algorithm has theoretical performance guarantees in contrast to prior work on concept drift which mainly deal with the problem in a ad-hoc manner. Indeed, it is customary to assume that online content is highly dynamic. A certain type of content may become popular for a certain period of time, and then, its popularity may decrease over time and a new content may emerge as popular. In addition, although the type of the content

remains the same, such as soccer news, its popularity may change over time due to exogenous events such as the World Cup etc. Similarly, a certain type of content may become popular for a certain type of demographics (e.g., users of a particular age, gender, profession, etc.). However, over time the interest of these users may shift to other types of content. In such cases, where the popularity of content changes over time for a user with context  $x$ ,  $\pi_c(x, t)$  denotes the probability that the user at time  $t$  will like content  $c$ .

As we stated earlier, a CA  $i$  can either recommend content from multimedia sources that it is directly connected to or can ask another CA for content. By asking for content  $c$  from another CA  $j$ , CA  $i$  will incur cost  $d_j^i \geq 0$ . For the purpose of our paper, the cost is a *generic* term. For instance, it can be a payment made to CA  $j$  to display it to CA  $i$ ’s user, or it may be associated with the advertising loss CA  $i$  incurs by directing its user to CA  $j$ ’s website. When the cost is payment, it can be money, tokens [27] or Bitcoins [28]. Since this cost is bounded, without loss of generality we assume that  $d_j^i \in [0, 1]$  for all  $i, j \in \mathcal{M}$ . In order make our model general, we also assume that there is a cost associated with recommending a type of content  $c \in \mathcal{C}_i$ , which is given by  $d_c^i \in [0, 1]$ , for CA  $i$ . For instance, this can be a payment made to the multimedia source that owns content  $c$ .

An intrinsic assumption we make is that the CAs are cooperative. That is, CA  $j \in \mathcal{M}_{-i}$  will return the content that is mostly to be liked by CA  $i$ ’s user when asked by CA  $i$  to recommend a content. This cooperative structure can be justified as follows. Whenever a user likes the content of CA  $j$  (either its own user or user of another CA), CA  $j$  obtains a benefit. This can be either an additional payment made by CA  $i$  when the content recommended by CA  $j$  is liked by CA  $i$ ’s user, or it can simply be the case that whenever a content of CA  $j$  is liked by someone its popularity increases. However, we assume that the CAs’ decisions do not change their pool of users. The future user arrivals to the CAs are independent of their past content matching strategies. For instance, users of a CA may have monthly or yearly subscriptions, so they will not shift from one CA to another CA when they like the content of the other CA.

The goal of CA  $i$  is to explore the matching actions in  $\mathcal{K}_i$  to learn the best content for each context, while at the same time exploiting the best content for the user with context  $x_i(t)$  arriving at each time instance  $t$  to maximize its total number of likes minus costs. CA  $i$ ’s problem can be modeled as a contextual bandit problem [10], [21], [22], [29], where likes and costs translate into rewards. In the next subsection, we formally define the benchmark solution which is computed using perfect knowledge about the probability that a content  $c$  will be liked by a user with context  $x$  (which requires complete knowledge of user and content characteristics). Then, we define the regret which is the performance loss due to uncertainty about the user and content characteristics.

#### B. Optimal Content Matching with Complete Information

Our benchmark when evaluating the performance of the learning algorithms is the optimal solution which always recommends the content with the highest relevance score minus

cost for CA  $i$  from the set  $\mathcal{C}$  given context  $x_i(t)$  at time  $t$ . This corresponds to selecting the best matching action in  $\mathcal{K}_i$  given  $x_i(t)$ . Next, we define the expected rewards of the matching actions, and the *action selection* policy of the benchmark. For a matching action  $k \in \mathcal{M}_{-i}$ , its relevance score is given as  $\pi_k(x) := \pi_{c_k^*(x)}(x)$ , where  $c_k^*(x) := \arg \max_{c \in \mathcal{C}_j} \pi_c(x)$ . For a matching action  $k \in \mathcal{C}_i$  its relevance score is equal to the relevance score of content  $k$ . The expected reward of CA  $i$  from choosing action  $k \in \mathcal{K}_i$  is given by the quasilinear utility function

$$\mu_k^i(x) := \pi_k(x) - d_k^i \quad (1)$$

where  $d_k^i \in [0, 1]$  is the normalized cost of choosing action  $k$  for CA  $i$ . Our proposed system will also work for more general expected reward functions as long as the expected reward of a learner is a function of the relevance score of the chosen action and the cost (payment, communication cost, etc.) associated with choosing that action. The *oracle* benchmark is given by

$$k_i^*(x) := \arg \max_{k \in \mathcal{K}_i} \mu_k^i(x) \quad \forall x \in \mathcal{X}. \quad (2)$$

The oracle benchmark depends on relevance scores as well as costs of matching content from its own content network or other CA's content network. The case  $d_k^i = 0$  for all  $k \in \mathcal{K}_i$  and  $i \in \mathcal{M}$ , corresponds to the scheme in which content matching has zero cost, hence  $k_i^*(x) = \arg \max_{k \in \mathcal{K}_i} \pi_k(x) = \arg \max_{c \in \mathcal{C}} \pi_c(x)$ . This corresponds to the best centralized solution, where CAs act as a single entity. On the other hand, when  $d_k^i \geq 1$  for all  $k \in \mathcal{M}_{-i}$  and  $i \in \mathcal{M}$ , in the oracle benchmark a CA must not cooperate with any other CA and should only use its own content. Hence  $k_i^*(x) = \arg \max_{c \in \mathcal{C}_i} (\pi_c(x) - d_c^i)$ . In the following subsections, we will show that independent of the values of relevance scores and costs, our algorithms will achieve sublinear regret (in the number of users or equivalently time) with respect to the oracle benchmark.

### C. The Regret of Learning

In this subsection we define the regret as a performance measure of the learning algorithm used by the CAs. Simply, the regret is the loss incurred due to the unknown system dynamics. Regret of a learning algorithm which selects the matching action/arm  $a_i(t)$  at time  $t$  for CA  $i$  is defined with respect to the best matching action  $k_i^*(x)$  given in (2). Then, the regret of CA  $i$  at time  $T$  is

$$R_i(T) := \sum_{t=1}^T \left( \pi_{k_i^*(x_i(t))}(x_i(t)) - d_{k_i^*(x_i(t))}^i \right) - \mathbb{E} \left[ \sum_{t=1}^T \left( \mathbb{I}(y_i(t) = L) - d_{a_i(t)}^i \right) \right]. \quad (3)$$

Regret gives the convergence rate of the total expected reward of the learning algorithm to the value of the optimal solution given in (2). Any algorithm whose regret is sublinear, i.e.,  $R_i(T) = O(T^\gamma)$  such that  $\gamma < 1$ , will converge to the optimal solution in terms of the average reward.

A summary of notations is given in Table II. In the next section, we propose an online learning algorithm which achieves

$\mathcal{M}$ : Set of all CAs
$\mathcal{C}_i$ : Contents in the Content Network of CA $i$
$C_{\max}$ : $\max_{i \in \mathcal{M}}  \mathcal{C}_i $
$\mathcal{C}$ : Set of all contents
$\mathcal{X} = [0, 1]^d$ : Context space
$\mathcal{Y}$ : Set of feedbacks a user can give
$x_i(t)$ : $d$ -dimensional context of $t$ th user of CA $i$
$y_i(t)$ : Feedback of the $t$ th user of CA $i$
$\mathcal{K}_i$ : Set of content matching actions of CA $i$
$\pi_c(x)$ : Relevance score of content $c$ for context $x$
$d_k^i$ : Cost of choosing matching action $k$ for CA $i$
$\mu_k^i(x)$ : Expected reward (static) of CA $i$ from matching action $k$ for context $x$
$k_i^*(x)$ : Optimal matching action of CA $i$ given context $x$ (oracle benchmark)
$R_i(T)$ : Regret of CA $i$ at time $T$
$\beta_a := \sum_{t=1}^{\infty} 1/t^a$

TABLE II  
NOTATIONS USED IN PROBLEM FORMULATION.

sublinear regret when the user and content characteristics are static.

## IV. A DISTRIBUTED ONLINE CONTENT MATCHING ALGORITHM

In this section we propose an online learning algorithm for content matching when the user and content characteristics are static. In contrast to prior online learning algorithms that exploit the context information [10], [11], [20]–[22], [29], which consider a single learner setting, the proposed algorithm helps a CA to learn from the *experience* of other CAs. With this mechanism, a CA is able to recommend content from multimedia sources that it has no direct connection, without needing to know the IDs of such multimedia sources and their content. It learns about these multimedia sources only through the other CAs that it is connected to.

In order to bound the regret of this algorithm analytically we use the following assumption. When the content characteristics are static, we assume that each type of content has similar relevance scores for similar contexts; we formalize this in terms of a Lipschitz condition.

*Assumption 1:* There exists  $L > 0$ ,  $\gamma > 0$  such that for all  $x, x' \in \mathcal{X}$  and  $c \in \mathcal{C}$ , we have  $|\pi_c(x) - \pi_c(x')| \leq L \|x - x'\|^\gamma$ .

Assumption 1 indicates that the probability that a type  $c$  content is liked by users with similar contexts will be similar to each other. For instance, if two users have similar age, gender, etc., then it is more likely that they like the same content. We call  $L$  the *similarity constant* and  $\gamma$  the *similarity exponent*. These parameters will depend on the characteristics of the users and the content. We assume that  $\gamma$  is known by the CAs. However, an unknown  $\gamma$  can be estimated online using the history of likes and dislikes by users with different contexts, and our proposed algorithms can be modified to include the estimation of  $\gamma$ .

In view of this assumption, the important question becomes how to learn from the past experience which content to match with the current user. We answer this question by proposing an algorithm which partitions the context space of a CA, and learns the relevance scores of different types of content for each set in the partition, based only on the past experience in that set. The algorithm is designed in a way to

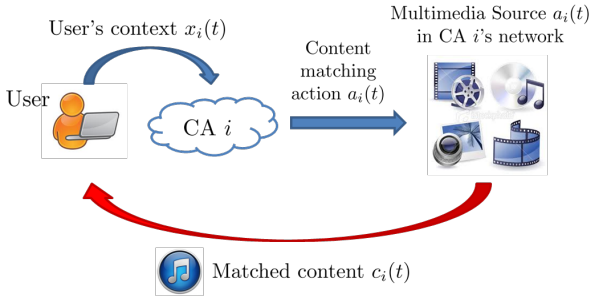


Fig. 2. Content matching within own content network.

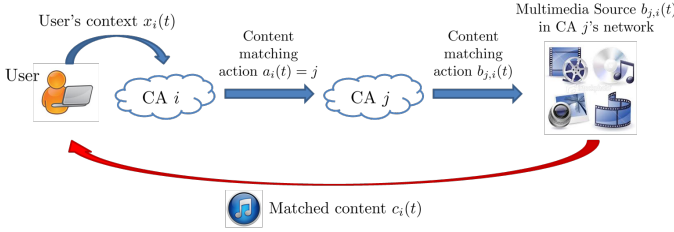


Fig. 3. Content matching from the network of another CA.

achieve optimal tradeoff between the size of the partition and the past observations that can be used together to learn the relevance scores. It also includes a mechanism to help CAs learn from each other's users. We call our proposed algorithm the *DIStributed Content Matching algorithm* (DISCOM), and its pseudocode is given in Fig. 4, Fig. 5 and Fig. 6.

Each CA  $i$  has two tasks: matching content with its own users and matching content with the users of other CAs when requested by those CAs. We call the first task the *maximization task* (implemented by  $\text{DISCOM}_{\max}$  given in Fig. 5), since the goal of CA  $i$  is to maximize the number of likes from its own users. The second task is called the *cooperation task* (implemented by  $\text{DISCOM}_{\text{coop}}$  given in Fig. 6), since the goal of CA  $i$  is to help other CAs obtain content from its own content network in order to maximize the likes they receive from their users. This cooperation is beneficial to CA  $i$  because of numerous reasons. Firstly, since every CA cooperates, CA  $i$  can reach a much larger set of content including the content from other CA's content networks, hence will be able to provide content with higher relevance score to its users. Secondly, when CA  $i$  helps CA  $j$ , it will observe the feedback of CA  $j$ 's user for the matched content, hence will be able to update the estimated relevance score of its content, which is beneficial if a user similar to CA  $j$ 's user arrives to CA  $i$  in the future. Thirdly, payment mechanisms can be incorporated to the system such that CA  $i$  gets a payment from CA  $j$  when its content is liked by CA  $j$ 's user.

In summary, there are two types of content matching actions for a user of CA  $i$ . In the first type, the content is recommended from a source that is directly connected to CA  $i$ , while in the second type, the content is recommended from a source that CA  $i$  is connected through another CA. The information exchange between multimedia sources and CAs for these two types of actions is shown in Fig. 2 and Fig. 3.

Let  $T$  be the time horizon of interest (equivalent to the number of users that arrive to each CA). DISCOM creates a partition of  $\mathcal{X} = [0, 1]^d$  based on  $T$ . For instance  $T$  can be

the average number of visits to the CA's website in one day. Although in reality the average number of visits to different CAs can be different, our analysis of the regret in this section will hold since it is the worst-case analysis (assuming that users arrive only to CA  $i$ , while the other CAs only learn through CA  $i$ 's users). Moreover, the case of heterogeneous number of visits can be easily addressed if each CA informs other CAs about its average number of visits. Then, CA  $i$  can keep  $M$  different partitions of the context space; one for itself and  $M - 1$  for the other CAs. If called by CA  $j$ , it will match a content to CA  $j$ 's user based on the partition it keeps for CA  $j$ . Hence, we focus on the case when  $T$  is common to all CAs.

We first define  $m_T$  as the *slicing level* used by DISCOM, which is an integer that is used to partition  $\mathcal{X}$ . DISCOM forms a partition of  $\mathcal{X}$  consisting of  $(m_T)^d$  sets (hypercubes) where each set is a  $d$ -dimensional hypercube with edge length  $1/m_T$ . This partition is denoted by  $\mathcal{P}_T$ . The hypercubes in  $\mathcal{P}_T$  are oriented such that one of them has a corner located at the origin of the  $d$ -dimensional Euclidian space. It is clear that the number of hypercubes is increasing in  $m_T$ , while their size is decreasing in  $m_T$ . When  $m_T$  is small each hypercube covers a large set of contexts, hence the number of past observations that can be used to estimate relevance scores of matching actions in each set is large. However, the variation of the true relevance scores of the contexts within a hypercube increases with the size of the hypercube. DISCOM should set  $m_T$  to a value that balances this tradeoff.

A hypercube in  $\mathcal{P}_T$  is denoted by  $p$ . The hypercube in  $\mathcal{P}_T$  that contains context  $x_i(t)$  is denoted by  $p_i(t)$ . When  $x_i(t)$  is located at a boundary of multiple hypercubes in  $\mathcal{P}_T$ , it is randomly assigned to one of these hypercubes.

$\text{DISCOM}_{\max}$  operates as follows. CA  $i$  matches its user at time  $t$  with a content by taking a matching action based on one of the three phases: *training* phase in which CA  $i$  requests content from another CA  $j$  for the purpose of helping CA  $j$  to learn the relevance score of content in its content network for users with context  $x_i(t)$  (but CA  $i$  does not update the relevance score for CA  $j$  because it thinks that CA  $j$  may not know much about its own content), *exploration* phase in which CA  $i$  selects a matching action in  $\mathcal{K}_i$  and updates its relevance score based on the feedback of its user, and *exploitation* phase in which CA  $i$  chooses the matching action with the highest relevance score minus cost.

Since the CAs are cooperative, when another CA requests content from CA  $i$ , CA  $i$  will choose content from its content network with the highest estimated relevance score for the user of the requesting CA. To maximize the number of likes minus costs in exploitations, CA  $i$  must have an accurate estimate of the relevance scores of other CAs. This task is not trivial since CA  $i$  does not know the content network of other CAs. In order to do this, CA  $i$  should smartly select which of its users' feedbacks to use when estimating the relevance score of CA  $j$ . The feedbacks should come from previous times at which CA  $i$  has a very high confidence that the content of CA  $j$  matched with its user is the one with the highest relevance score for the context of CA  $i$ 's user. Thus, the training phase of CA  $i$  helps other CAs build accurate estimates about the relevance scores



DISCOM for CA  $i$ :

- 1: Input:  $H_1(t), H_2(t), H_3(t), T, m_T$
- 2: Initialize: Partition  $\mathcal{X}$  into hypercubes denoted by  $\mathcal{P}_T$
- 3: Initialize: Set counters  $N_p^i = 0, \forall p \in \mathcal{P}_T$ ,  
 $N_{k,p}^i = 0, \forall k \in \mathcal{K}_i, p \in \mathcal{P}_T, N_{j,p}^{\text{tr},i} = 0, \forall j \in \mathcal{M}_{-i}, p \in \mathcal{P}_T$
- 4: Initialize: Set relevance score estimates  $\bar{r}_{k,p}^i = 0, \forall k \in \mathcal{K}_i, p \in \mathcal{P}_T$
- 5: **while**  $t \geq 1$  **do**
- 6: Run DISCOM<sub>max</sub> to find  $p = p_i(t)$ , to obtain a matching action  $a_i$ , and value of *train* flag
- 7: If  $a_i \in \mathcal{M}_{-i}$  ask CA  $a_i$  for content and pass  $x_i(t)$
- 8: Receive  $\mathcal{CA}_i(t)$ , the set of CAs who requested content from CA  $i$ , and their contexts
- 9: **if**  $\mathcal{CA}_i(t) \neq \emptyset$  **then**
- 10: Run DISCOM<sub>coop</sub> to obtain the content to be selected  $b_i := \{b_{i,j}\}_{j \in \mathcal{CA}_i(t)}$  and hypercubes that the contexts of the users in  $\mathcal{CA}_i(t)$  lie in  $\mathbf{p}_i := \{p_{i,j}\}_{j \in \mathcal{CA}_i(t)}$
- 11: **end if**
- 12: **if**  $a_i \in \mathcal{C}_i$  **then**
- 13: Pay cost  $d_{a_i}^i$ , obtain content  $a_i$
- 14: Show  $a_i$  to the user, receive feedback  $r \in \{0, 1\}$  drawn from  $\text{Ber}_{a_i}(x_i(t))^a$
- 15: **else**
- 16: Pay cost  $d_{a_i}^i$ , obtain content  $b_{a_i,i}$  from CA  $a_i$
- 17: Show  $b_{a_i,i}$  to the user, receive feedback  $r \in \{0, 1\}$  drawn from  $\text{Ber}_{b_{a_i,i}}(x_i(t))$
- 18: **end if**
- 19: **if** *train* = 1 **then**
- 20:  $N_{a_i,p}^{\text{tr},i} ++$
- 21: **else**
- 22:  $\bar{r}_{a_i,p}^i = (\bar{r}_{a_i,p}^i N_{a_i,p}^i + r) / (N_{a_i,p}^i + 1)$
- 23:  $N_p^i ++, N_{a_i,p}^i ++$
- 24: **end if**
- 25: **if**  $\mathcal{CA}_i(t) \neq \emptyset$  **then**
- 26: **for**  $j \in \mathcal{CA}_i(t)$  **do**
- 27: Send content  $b_{i,j}$  to CA  $j$ 's user
- 28: Observe feedback  $r$  drawn from  $\text{Ber}_{b_{i,j}}(x_j(t))$
- 29:  $\bar{r}_{b_{i,j},p_{i,j}}^i = \frac{\bar{r}_{b_{i,j},p_{i,j}}^i N_{b_{i,j},p_{i,j}}^i + r}{N_{b_{i,j},p_{i,j}}^i + 1}$
- 30:  $N_{p_{i,j}}^i ++, N_{b_{i,j},p_{i,j}}^i ++$
- 31: **end for**
- 32: **end if**
- 33:  $t = t + 1$
- 34: **end while**

<sup>a</sup> $\text{Ber}_{a_i}(x_i(t))$  is the Bernoulli distribution with expected value  $\pi_{a_i}(x_i(t))$

Fig. 4. Pseudocode for DISCOM algorithm.

of their content, before CA  $i$  uses any feedback for content coming from these CAs to form relevance score estimates about them. In contrast, the exploration phase of CA  $i$  helps it to build accurate estimates about the relevance score of its matching actions.

At time  $t$ , the phase that CA  $i$  will be in is determined by the amount of time it had explored, exploited or trained for past users with contexts similar to the context of the current user. For this CA  $i$  keeps counters and control functions which are described below. Let  $N_p^i(t)$  be the number of user arrivals to CA  $i$  with contexts in  $p \in \mathcal{P}_T$  by time  $t$  (its own arrivals and arrivals to other CAs who requested content from CA  $i$ ) except the training phases of CA  $i$ . For  $c \in \mathcal{C}_i$ , let  $N_{c,p}^i(t)$  be the number of times content  $c$  is selected in response to a user arriving to CA  $i$  with context in hypercube  $p$  by time  $t$  (including times other CAs request content from CA  $i$  for their users with contexts in set  $p$ ). Other than these, CA  $i$  keeps two counters for each other CA in each set in the partition, which it uses to decide the phase it should be in. The first one, i.e.,

DISCOM<sub>max</sub> (maximization part of DISCOM) for CA  $i$ :

- 1: *train* = 0
- 2: Find the hypercube in  $\mathcal{P}_T$  that  $x_i(t)$  belongs to, i.e.,  $p_i(t)$
- 3: Let  $p = p_i(t)$
- 4: Compute the set of under-explored matching actions  $\mathcal{C}_{i,p}^{\text{ue}}(t)$  given in (4)
- 5: **if**  $\mathcal{C}_{i,p}^{\text{ue}}(t) \neq \emptyset$  **then**
- 6: Select  $a_i$  randomly from  $\mathcal{C}_{i,p}^{\text{ue}}(t)$
- 7: **else**
- 8: Compute the set of training candidates  $\mathcal{M}_{i,p}^{\text{ct}}(t)$  given in (5)
- 9: //Update the counters of training candidates
- 10: **for**  $j \in \mathcal{M}_{i,p}^{\text{ct}}(t)$  **do**
- 11: Obtain  $N_p^j$  from CA  $j$ , set  $N_{j,p}^{\text{tr},i} = N_p^j - N_{j,p}^i$
- 12: **end for**
- 13: Compute the set of under-trained CAs  $\mathcal{M}_{i,p}^{\text{ut}}(t)$  given in (6)
- 14: Compute the set of under-explored CAs  $\mathcal{M}_{i,p}^{\text{ue}}(t)$  given in (7)
- 15: **if**  $\mathcal{M}_{i,p}^{\text{ut}}(t) \neq \emptyset$  **then**
- 16: Select  $a_i$  randomly from  $\mathcal{M}_{i,p}^{\text{ut}}(t)$ , *train* = 1
- 17: **else if**  $\mathcal{M}_{i,p}^{\text{ue}}(t) \neq \emptyset$  **then**
- 18: Select  $a_i$  randomly from  $\mathcal{M}_{i,p}^{\text{ue}}(t)$
- 19: **else**
- 20: Select  $a_i$  randomly from  $\arg \max_{k \in \mathcal{K}_i} \bar{r}_{k,p}^i - d_k^i$
- 21: **end if**
- 22: **end if**

Fig. 5. Pseudocode for the maximization part of DISCOM algorithm.

DISCOM<sub>coop</sub> (cooperation part of DISCOM) for CA  $i$

- 1: **for**  $j \in \mathcal{CA}_i(t)$  **do**
- 2: Find the set in  $\mathcal{P}_T$  that  $x_j(t)$  belongs to, i.e.,  $p_{i,j}$
- 3: Compute the set of under-explored matching actions  $\mathcal{C}_{i,p_{i,j}}^{\text{ue}}(t)$  given in (4)
- 4: **if**  $\mathcal{C}_{i,p_{i,j}}^{\text{ue}}(t) \neq \emptyset$  **then**
- 5: Select  $b_{i,j}$  randomly from  $\mathcal{C}_{i,p_{i,j}}^{\text{ue}}(t)$
- 6: **else**
- 7:  $b_{i,j} = \arg \max_{c \in \mathcal{C}_i} \bar{r}_{c,p_{i,j}}^i$
- 8: **end if**
- 9: **end for**

Fig. 6. Pseudocode for the cooperation part of DISCOM algorithm.

$N_{j,p}^{\text{tr},i}(t)$ , is an estimate on the number of user arrivals with contexts in  $p$  to CA  $j$  from all CAs except the training phases of CA  $j$  and exploration, exploitation phases of CA  $i$ . This counter is only updated when CA  $i$  thinks that CA  $j$  should be trained. The second one, i.e.,  $N_{j,p}^i(t)$ , counts the number of users of CA  $i$  with contexts in  $p$  for which content is requested from CA  $j$  at exploration and exploitation phases of CA  $i$  by time  $t$ .

At each time slot  $t$ , CA  $i$  first identifies  $p_i(t)$ . Then, it chooses its phase at time  $t$  by giving highest priority to exploration of content in its own content network, second highest priority to training of the other CAs, third highest priority to exploration of the other CAs, and lowest priority to exploitation. The reason that exploration of own content has a higher priority than training of other CAs is that it will minimize the number of times CA  $i$  will be trained by other CAs, which we describe below.

First, CA  $i$  identifies the set of under-explored content in its content network:

$$\mathcal{C}_{i,p}^{\text{ue}}(t) := \{c \in \mathcal{C}_i : N_{c,p}^i(t) \leq H_1(t)\} \quad (4)$$

where  $H_1(t)$  is a deterministic, increasing function of  $t$  which is called *the control function*. The value of this function will affect the regret of DISCOM. For  $c \in \mathcal{C}_i$ , the accuracy of

relevance score estimates increase with  $H_1(t)$ , hence it should be selected to balance the tradeoff between accuracy and the number of explorations. If  $\mathcal{C}_{i,p}^{\text{ue}}(t)$  is non-empty, CA  $i$  enters the exploration phase and randomly selects a content in this set to explore. Otherwise, it identifies the set of training candidates:

$$\mathcal{M}_{i,p}^{\text{ct}}(t) := \{j \in \mathcal{M}_{-i} : N_{j,p}^{\text{tr},i}(t) \leq H_2(t)\} \quad (5)$$

where  $H_2(t)$  is a control function similar to  $H_1(t)$ . Accuracy of other CA's relevance score estimates of content in their own networks increases with  $H_2(t)$ , hence it should be selected to balance the possible reward gain of CA  $i$  due to this increase with the reward loss of CA  $i$  due to the number of trainings. If this set is non-empty, CA  $i$  asks the CAs  $j \in \mathcal{M}_{i,p}^{\text{ct}}(t)$  to report  $N_p^j(t)$ . Based in the reported values it recomputes  $N_{j,p}^{\text{tr},i}(t)$  as  $N_{j,p}^{\text{tr},i}(t) = N_p^j(t) - N_{j,p}^i(t)$ . Using the updated values, CA  $i$  identifies the set of under-trained CAs:

$$\mathcal{M}_{i,p}^{\text{ut}}(t) := \{j \in \mathcal{M}_{-i} : N_{j,p}^{\text{tr},i}(t) \leq H_2(t)\}. \quad (6)$$

If this set is non-empty, CA  $i$  enters the training phase and randomly selects a CA in this set to train it. When  $\mathcal{M}_{i,p}^{\text{ct}}(t)$  or  $\mathcal{M}_{i,p}^{\text{ut}}(t)$  is empty, this implies that there is no under-trained CA, hence CA  $i$  checks if there is an under-explored matching action. The set of CAs for which CA  $i$  does not have accurate relevance scores is given by

$$\mathcal{M}_{i,p}^{\text{ue}}(t) := \{j \in \mathcal{M}_{-i} : N_{j,p}^i(t) \leq H_3(t)\} \quad (7)$$

where  $H_3(t)$  is also a control function similar to  $H_1(t)$ . If this set is non-empty, CA  $i$  enters the exploration phase and randomly selects a CA in this set to request content from to explore it. Otherwise, CA  $i$  enters the exploitation phase in which it selects the matching action with the highest estimated relevance score minus cost for its user with context  $x_i(t) \in p = p_i(t)$ , i.e.,

$$a_i(t) \in \arg \max_{k \in \mathcal{K}_i} \bar{r}_{k,p}^i(t) - d_k^i \quad (8)$$

where  $\bar{r}_{k,p}^i(t)$  is the sample mean estimate of the relevance score of CA  $i$  for matching action  $k$  at time  $t$ , which is computed as follows. For  $j \in \mathcal{M}_{-i}$ , let  $\mathcal{E}_{j,p}^i(t)$  be the set of feedbacks collected by CA  $i$  at times it selected CA  $j$  while CA  $i$ 's users' contexts are in set  $p$  in its exploration and exploitation phases by time  $t$ . For estimating the relevance score of contents in its own content network, CA  $i$  can also use the feedback obtained from other CAs' users at times they requested content from CA  $i$ . In order to take this into account, for  $c \in \mathcal{C}_i$ , let  $\mathcal{E}_{c,p}^i(t)$  be the set of feedbacks observed by CA  $i$  at times it selected its content  $c$  for its own users with contexts in set  $p$  union the set of feedbacks observed by CA  $i$  when it selected its content  $c$  for the users of other CAs with contexts in set  $p$  who requests content from CA  $i$  by time  $t$ .

Therefore, sample mean relevance score of matching action  $k \in \mathcal{K}_i$  for users with contexts in set  $p$  for CA  $i$  is defined as  $\bar{r}_{k,p}^i(t) = \left( \sum_{r \in \mathcal{E}_{k,p}^i(t)} r \right) / |\mathcal{E}_{k,p}^i(t)|$ . An important observation is that computation of  $\bar{r}_{k,p}^i(t)$  does not take into account the matching costs. Let  $\hat{\mu}_{k,p}^i(t) := \bar{r}_{k,p}^i(t) - d_k^i$  be the estimated *net reward* (relevance score minus cost) of matching action  $k$  for

$L$ : Similarity constant. $\gamma$ : Similarity exponent
$T$ : Time horizon
$m_T$ : Slicing level of DISCOM
$\mathcal{P}_T$ : DISCOM's partition of $\mathcal{X}$ into $(m_T)^d$ hypercubes
$p_i(t)$ : Hypercube in $\mathcal{P}_T$ that contains $x_i(t)$
$N_p^i(t)$ : Number of <i>all</i> user arrivals to CA $i$ with contexts in $p \in \mathcal{P}_T$ by time $t$ except the training phases of CA $i$
$N_{c,p}^i(t)$ : Number of times content $c$ is selected in response to a user arriving to CA $i$ with context in hypercube $p$ by time $t$
$N_{j,p}^{\text{tr},i}(t)$ : Estimate of CA $i$ on the number of user arrivals with contexts in $p$ to CA $j$ from all CAs except the training phases of CA $j$ and exploration, exploitation phases of CA $i$
$N_{j,p}^i(t)$ : Number of users of CA $i$ with contexts in $p$ for which content is requested from CA $j$ at exploration and exploitation phases of CA $i$ by time $t$
$H_1(t), H_2(t), H_3(t)$ : Control functions of DISCOM
$\mathcal{C}_{i,p}^{\text{ue}}(t)$ : Set of under-explored content in $\mathcal{C}_i$
$\mathcal{M}_{i,p}^{\text{ct}}(t)$ : Set of training candidates of CA $i$
$\mathcal{M}_{i,p}^{\text{ut}}(t)$ : Set of CAs under-trained by CA $i$
$\mathcal{M}_{i,p}^{\text{ue}}(t)$ : Set of CAs under-explored by CA $i$
$\bar{r}_{k,p}^i(t)$ : Sample mean relevance score of action $k$ of CA $i$ at time $t$
$\hat{\mu}_{k,p}^i(t)$ : Estimated net reward of action $k$ of CA $i$ at time $t$

TABLE III  
NOTATIONS USED IN DEFINITION OF DISCOM.

set  $p$ . Of note, when there is more than one maximizer of (8), one of them is randomly selected. In order to run DISCOM, CA  $i$  does not need to keep the sets  $\mathcal{E}_{k,p}^i(t)$  in its memory.  $\bar{r}_{k,p}^i(t)$  can be computed by using only  $\bar{r}_{k,p}^i(t-1)$  and the feedback at time  $t$ .

The cooperation part of DISCOM, i.e., DISCOM<sub>coop</sub> operates as follows. Let  $\mathcal{CA}_i(t)$  be the set CAs who request content from CA  $i$  at time  $t$ . For each  $j \in \mathcal{CA}_i(t)$ , CA  $i$  first checks if it has any under-explored content  $c$  for  $p_j(t)$ , i.e.,  $c$  such that  $N_{c,p_j(t)}^i(t) \leq H_1(t)$ . If so, it randomly selects one of its under-explored content to match it with the user of CA  $j$ . Otherwise, it exploits its content in  $\mathcal{C}_i$  with the highest estimated relevance score for CA  $j$ 's current user's context, i.e.,

$$b_{i,j}(t) \in \arg \max_{c \in \mathcal{C}_i} \bar{r}_{c,p_j(t)}^i(t). \quad (9)$$

A summary of notations used in the description of DISCOM is given in Table III. The following theorem provides a bound on the regret of DISCOM.

*Theorem 1:* When DISCOM is run by all CAs with parameters  $H_1(t) = t^{2\gamma/(3\gamma+d)} \log t$ ,  $H_2(t) = C_{\max} t^{2\gamma/(3\gamma+d)} \log t$ ,  $H_3(t) = t^{2\gamma/(3\gamma+d)} \log t$  and  $m_T = \lceil T^{1/(3\gamma+d)} \rceil$ ,<sup>6</sup> we have

$$\begin{aligned} R_i(T) &\leq 4(M + C_{\max} + 1)\beta_2 \\ &+ T^{\frac{2\gamma+d}{3\gamma+d}} \left( \frac{14Ld^{\gamma/2} + 12 + 4(|\mathcal{C}_i| + M)MC_{\max}\beta_2}{(2\gamma + d)/(3\gamma + d)} \right) \\ &+ 2^{d+1} Z_i \log T \\ &+ T^{\frac{d}{3\gamma+d}} 2^{d+1} (|\mathcal{C}_i| + 2(M - 1)), \end{aligned}$$

i.e.,  $R_i(T) = \tilde{O} \left( MC_{\max} T^{\frac{2\gamma+d}{3\gamma+d}} \right)$ ,<sup>7</sup> where  $Z_i = |\mathcal{C}_i| + (M - 1)(C_{\max} + 1)$ .

<sup>6</sup>For a number  $r \in \mathbb{R}$ , let  $\lceil r \rceil$  be the smallest integer that is greater than or equal to  $r$ .

<sup>7</sup> $\tilde{O}(\cdot)$  is the Big-O notation in which the terms with logarithmic growth rates are hidden.



*Proof:* The proof is given Appendix B. ■

For any  $d > 0$  and  $\gamma > 0$ , the regret given in Theorem 1 is sublinear in time (or number of user arrivals). This guarantees that the regret per-user, i.e., the time averaged regret, converges to 0 ( $\lim_{T \rightarrow \infty} \mathbb{E}[R_i(T)]/T = 0$ ). It is also observed that the regret increases in the dimension  $d$  of the context. By Assumption 1, a context is similar to another context if they are similar in each dimension, hence number of hypercubes in the partition  $\mathcal{P}_T$  increases with  $d$ .

In our analysis of the regret of DISCOM we assumed that  $T$  is fixed and given as an input to the algorithm. DISCOM can be made to run independently of the final time  $T$  by using a standard method called *the doubling trick* (see, e.g., [10]). The idea is to divide time into rounds with geometrically increasing lengths and run a new instance of DISCOM at each round. For instance, consider rounds  $\tau \in \{1, 2, \dots\}$ , where each round has length  $2^\tau$ . Run a new instance of DISCOM at the beginning of each round with time parameter  $2^\tau$ . This modified version will also have  $\tilde{O}(T^{(2\gamma+d)/(3\gamma+d)})$  regret.

Maximizing the satisfaction of an individual user is as important as maximizing the overall satisfaction of all users. The next corollary shows that by using DISCOM, CAs guarantee that their users will almost always be provided with the best content available within the entire content network.

*Corollary 1:* Assume that DISCOM is run with the set of parameters given in Theorem 1. When DISCOM is in exploitation phase for CA  $i$ , we have

$$\begin{aligned} P(\mu_{a_i(t)}^i(x_i(t)) < \mu_{k_i^*(x_i(t))}^i(x_i(t)) - \delta_T) \\ \leq \frac{2|\mathcal{K}_i|}{t^2} + \frac{2|\mathcal{K}_i|MC_{\max}\beta_2}{t^{\gamma/(3\gamma+d)}}. \end{aligned}$$

where  $\delta_T = (6Ld^{\gamma/2} + 6)T^{-\gamma/(3\gamma+d)}$ .

*Proof:* The proof is given Appendix C. ■

**Remark 1: (Differential Services)** Maximizing the performance for an individual user is particularly important for providing *differential services* based on the types of the users. For instance, a CA may want to provide higher quality recommendations to a subscriber (high type user) who has paid for the subscription compared to a non-subscriber (low type user). To do this, the CA can exploit the best content for the subscribed user, while perform exploration on a different user that is not subscribed.

## V. REGRET WHEN FEEDBACK IS MISSING

When analyzing the performance of DISCOM, we assumed that the users always provide a feedback: like or dislike. However, in most of the online content aggregation platforms user feedback is not always available. In this section we consider the effect of missing feedback on the performance of the proposed algorithm. We assume that each user gives a feedback with probability  $p_r$  (which is unknown to the CAs). If the user at time  $t$  does not give feedback, we assume that DISCOM does not update its counters. This will result in a larger number of trainings and explorations compared to the case when feedback is always available. The following theorem gives an upper bound on the regret of DISCOM for this case.

**Theorem 2:** Let the DISCOM algorithm run with parameters  $H_1(t) = t^{2\gamma/(3\gamma+d)} \log t$ ,  $H_2(t) = C_{\max} t^{2\gamma/(3\gamma+d)} \log t$ ,  $H_3(t) = t^{2\gamma/(3\gamma+d)} \log t$ , and  $m_T = \lceil T^{1/(3\gamma+d)} \rceil$ . Then, if a user reveals its feedback with probability  $p_r$ , we have for CA  $i$ ,

$$\begin{aligned} R_i(T) &\leq 4(M + C_{\max} + 1)\beta_2 \\ &+ T^{\frac{2\gamma+d}{3\gamma+d}} \left( \frac{14Ld^{\gamma/2} + 12 + 4(|C_i| + M)MC_{\max}\beta_2}{(2\gamma + d)/(3\gamma + d)} \right. \\ &\quad \left. + \frac{2^{d+1}Z_i}{p_r} \log T \right) \\ &+ T^{\frac{d}{3\gamma+d}} 2^{d+1} \frac{|C_i| + 2(M - 1)}{p_r}, \end{aligned}$$

i.e.,  $R_i(T) = \tilde{O}\left(MC_{\max}T^{\frac{2\gamma+d}{3\gamma+d}}/p_r\right)$ , where  $Z_i = |C_i| + (M - 1)(C_{\max} + 1)$ ,  $\beta_a := \sum_{t=1}^{\infty} 1/t^a$ .

*Proof:* The proof is given Appendix D. ■

From Theorem 2, we see that missing feedback does not change the time order of the regret. However, the regret is scaled with  $1/p_r$ , which is the expected number of users required for a single feedback.

## VI. LEARNING UNDER DYNAMIC USER AND CONTENT CHARACTERISTICS

When the user and content characteristics change over time, the relevance score of content  $c$  for a user with context  $x$  changes over time. In this section, we assume that the following relation holds between the probabilities that a content will be liked with users with similar contexts at two different times  $t$  and  $t'$ .

**Assumption 2:** For each  $c \in \mathcal{C}$ , there exists  $L > 0$ ,  $\gamma > 0$  such that for all  $x, x' \in \mathcal{X}$ , we have

$$|\pi_{c,t}(x) - \pi_{c,t'}(x')| \leq L(\|x - x'\|)^\gamma + |t - t'|/T_s$$

where  $1/T_s > 0$  is the speed of the change in user and content characteristics. We call  $T_s$  the *stability parameter*.

Assumption 2 captures the temporal dynamics of content matching which is absent in Assumption 1. Such temporal variations are often referred to as *concept drift* [30], [31]. When there is concept drift, a learner should also consider which past information to take into account when learning, in addition to how to combine the past information to learn the best matching strategy.

The following modification of DISCOM will deal with dynamically changing user and content characteristics by using a time window of past observations in estimating the relevance scores. The modified algorithm is called DISCOM with time window (DISCOM-W). This algorithm groups the time slots into rounds  $\zeta = 1, 2, \dots$  each having a fixed length of  $2\tau_h$  time slots, where  $\tau_h$  is an integer called *the half window length*. Some of the time slots in these rounds overlap with each other as given in Fig. 7. The idea is to keep separate control functions and counters for each round, and calculate the sample mean relevance scores for groups of similar contexts based only on the observations that are made during the time window of that round. We call  $\eta = 1$  the *initialization round*. The control functions for the initialization round of DISCOM-W is the same as the control functions  $H_1(t)$ ,  $H_2(t)$  and  $H_3(t)$

of DISCOM whose values are given in Theorem 1. For the other rounds  $\zeta > 1$ , the control functions depend on  $\tau_h$  and are given as

$$H_1^{\tau_h}(t) = H_3^{\tau_h}(t) = (t \bmod \tau_h + 1)^z \log(t \bmod \tau_h + 1)$$

and

$$H_2^{\tau_h}(t) = C_{\max}(t \bmod \tau_h + 1)^z \log(t \bmod \tau_h + 1)$$

for some  $0 < z < 1$ . Each round  $\eta$  is divided into two sub-rounds. Except the initialization round, i.e.,  $\eta = 1$ , the first sub-round is called the *passive* sub-round, while the second sub-round is called the *active* sub-round. For the initialization round both sub-rounds are active sub-rounds. In order to reduce the number of trainings and explorations, DISCOM-W has an overlapping round structure as shown in Fig. 7. For each round except the initialization round, passive sub-rounds of round  $\eta$ , overlaps with the active sub-round of round  $\eta - 1$ . DISCOM-W operates in the same way as DISCOM in each round. DISCOM-W can be viewed as an algorithm which generates a new instance of DISCOM at the beginning of each round, with the modified control functions. DISCOM-W runs two different instances of DISCOM at each round. One of these instances is the active instance based on which content matchings are performed, and the other one is the passive instance which learns through the content matchings made by the active instance.

Let the instance of DISCOM that is run by DISCOM-W at round  $\eta$  be  $\text{DISCOM}_\eta$ . The hypercubes of  $\text{DISCOM}_\eta$  are formed in a way similar to DISCOM's. The input time horizon is taken as  $T_s$  which is the stability parameter given in Assumption 2, and the slicing parameter  $m_{T_s}$  is set accordingly.  $\text{DISCOM}_\eta$  uses the partition of  $\mathcal{X}$  into  $(m_{T_s})^d$  hypercubes denoted by  $\mathcal{P}_{T_s}$ . When all CAs are using DISCOM-W, the matching action selection of CA  $i$  only depends on the history of content matchings and feedback observations at round  $\eta$ . If time  $t$  is in the active sub-round of round  $\eta$ , matching action of CA  $i \in \mathcal{M}$  is taken according to  $\text{DISCOM}_\eta$ . As a result of the content matching, sample mean relevance scores and counters of both  $\text{DISCOM}_\eta$  and  $\text{DISCOM}_{\eta+1}$  are updated. Else if time  $t$  is in the passive sub-round of round  $\eta$ , matching action of CA  $i \in \mathcal{M}$  is taken according to  $\text{DISCOM}_{\eta-1}$  (see Fig. 7). As a result of this, sample mean relevance scores and counters of both  $\text{DISCOM}_{\eta-1}$  and  $\text{DISCOM}_\eta$  are updated.

At the start of a round  $\eta$ , the relevance score estimates and counters for  $\text{DISCOM}_\eta$  are equal to zero. However, due to the two sub-round structure, when the active sub-round of round  $\eta$  starts, CA  $i$  already has some observations for the context and actions taken in the passive sub-round of that round, hence depending on the arrivals and actions in the passive sub-round, the CA may even start the active sub-round by exploiting, whereas it should have always spent some time in training and exploration if it starts an active sub-round without any past observations (cold start problem).

In this section, due to the concept drift, even though the context of a past user can be similar to the context of the current user, their relevance scores for a content  $c$  can be very different. Hence DISCOM-W assumes that a past user

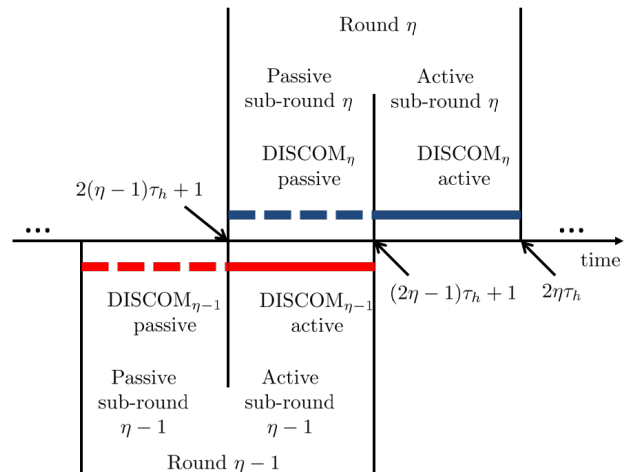


Fig. 7. Operation of DISCOM-W showing the round structure and the different instances of DISCOM running for each round.

is similar to the current user only if it arrived in the current round. Since round length is fixed, it is impossible to have sublinear number of similar context observations for every  $t$ . Thus, achieving sublinear regret under concept drift is not possible. Therefore, in this section we focus on the average regret which is given by

$$R_i^{\text{avg}}(T) := \frac{1}{T} \sum_{t=1}^T \left( \pi_{k_i^*(x_i(t))}(x_i(t)) - d_{k_i^*(x_i(t))}^i \right) - \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \left( \mathbb{I}(y_i(t) = L) - d_{a_i(t)}^i \right) \right].$$

The following theorem bounds the average regret of DISCOM-W.

**Theorem 3:** When DISCOM-W is run with parameters

$$H_1^{\tau_h}(t) = H_3^{\tau_h}(t) = (t \bmod \tau_h + 1)^{\frac{2\gamma}{3\gamma+d}} \log(t \bmod \tau_h + 1)$$

$$H_2^{\tau_h}(t) = C_{\max}(t \bmod \tau_h + 1)^{\frac{2\gamma}{3\gamma+d}} \log(t \bmod \tau_h + 1)$$

$m_{T_s} = \lceil T_s^{\frac{1}{3\gamma+d}} \rceil$  and  $\tau_h = \lceil T_s^{(3\gamma+d)/(4\gamma+d)} \rceil$ ,<sup>8</sup> where  $T_s$  is the stability parameter which is given in Assumption 2, the time averaged regret of CA  $i$  by time  $T$  is

$$R_i^{\text{avg}}(T) = \tilde{O} \left( T_s^{-\frac{\gamma}{4\gamma+d}} \right)$$

for any  $T > 0$ . Hence DISCOM-W is  $\epsilon = \tilde{O} \left( T_s^{-\frac{\gamma}{4\gamma+d}} \right)$  approximately optimal in terms of the average reward.

*Proof:* The proof is given Appendix E.  $\blacksquare$

From the result of this theorem we see that the average regret decays as the stability parameter  $T_s$  increases. This is because, DISCOM-W will use a longer time window (round) when  $T_s$  is large, and thus can get more observations to estimate the sample mean relevance scores of the matching actions in that round, which will result in better estimates hence smaller number of suboptimal matching action selections. Moreover, the average number of trainings and explorations required decrease with the round length.

<sup>8</sup>For a number  $b$ ,  $\lfloor b \rfloor$  denotes the largest integer that is smaller than or equal to  $b$ .

## VII. NUMERICAL RESULTS

In this section we provide numerical results for our proposed algorithms DISCOM and DISCOM-W on real-world datasets.

### A. Datasets

For all the datasets below, for a CA the cost of choosing a content within the content network and the cost of choosing another CA is set to 0. Hence, the only factor that affects the total reward is the users' ratings for the contents.

**Yahoo! Today Module (YTM)** [5]: The dataset contains news article webpage recommendations of Yahoo! Front Page. Each instance is composed of (i) ID of the recommended content, (ii) the user's context (2-dimensional vector), (iii) the user's click information. The user's click information for a webpage/content is associated with the relevance score of that content. It is equal to 1 if the user clicked on the recommended webpage and 0 else. The dataset contains  $T = 70000$  instances and 40 different types of content. We generate 4 CAs and assign 10 of the 40 types of content to each CA's content network. Each CA has direct access to content in its own network, while it can also access to the content in other CAs' content network by requesting content from these CAs. Users are divided into four groups according to their contexts and each group is randomly assigned to one of the CAs. Hence, the user arrival processes to different CA's are different. The performance of a CA is evaluated in terms of the average number of clicks, i.e., click through rate (CTR), of the contents that are matched with its users.

**Music Dataset (MD)**: The dataset contains contextual information and ratings (like/dislike) of music genres (classical, rock, pop, rap) collected from 413 students at UCLA. We generate 2 CAs each specialized in two of the four music genres. Users among the 413 users randomly arrive to each CA. A CA either recommends a music content that is in its content network or asks another CA, specialized in another music genre, to provide a music item. As a result, the rating of the user for the genre of the provided music content is revealed to the CA. The performance of a CA is evaluated in terms of the average number of likes it gets for the contents that are matched with its users.

**Yahoo! Today Module (YTM) with Drift (YTMD)**: This dataset is generated from YTM to simulate the scenario where the user ratings for a particular content changes over time. After every 10000 instances, 20 contents are randomly selected and user clicks for these contents are set to 0 (no click) for the next 10000 instances. For instance, this can represent a scenario where some news articles lose their popularity a day after they become available while some other news articles related to ongoing events will stay popular for several days.

### B. Learning Algorithms

While DISCOM and DISCOM-W are the first distributed algorithms to perform content aggregation (see Table I), we compare their performance with distributed versions of the centralized algorithms proposed in [5], [10], [19], [22]. In the distributed implementation of these centralized algorithms, we assume that each CA runs an independent instance of these

algorithms. For instance, when implementing a centralized algorithm  $A$  on the distributed system of CAs, we assume that each CA  $i$  runs its own instance of  $A$  denoted by  $A_i$ . When CA  $i$  selects CA  $j$  as a matching action in  $\mathcal{K}_i$  by using its algorithm  $A_i$ , CA  $j$  will select the content for CA  $i$  using its algorithm  $A_j$  with CA  $i$ 's user's context on the set of contents  $\mathcal{C}_j$ . In our numerical results, each algorithm is run for different values of its input parameters. The results are shown for the parameter values for which the corresponding algorithm performs the best.

**DISCOM**: Our algorithm given in Fig. 4 with control functions  $H_1(t)$ ,  $H_2(t)$  and  $H_3(t)$  divided by 10 for MD, and by 20 for YTM and YTMD to reduce the number of trainings and explorations.<sup>9</sup>

**DISCOM-W**: Our algorithm given in Fig. 7 which is the time-windowed version of DISCOM with control functions  $H_1(t)$ ,  $H_2(t)$  and  $H_3(t)$  divided by 20 to reduce the number of trainings and explorations.

As we mentioned in Remark 1, both DISCOM and DISCOM-W can provide differential services to its users. In this case both algorithms always exploit for the users with high type (subscribers) and if necessary can train and explore for the users with low type (non-subscribers). Hence, the performance of DISCOM and DISCOM-W for differential services is equal to their performance for the set of high type users.

**LinUCB** [5], [22]: This algorithm computes an index for each matching action by assuming that the relevance score of a matching action for a user is a linear combination of the contexts of the user. Then for each user it selects the matching action with the highest index.

**Hybrid- $\epsilon$**  [19]: This algorithm forms context-dependent sample mean rewards for the matching actions by considering the history of observations and decisions for groups of contexts that are similar to each other. For user  $t$  it either explores a random matching action with probability  $\epsilon_t$  or exploits the best matching action with probability  $1 - \epsilon_t$ , where  $\epsilon_t$  is decreasing in  $t$ .

**Contextual zooming (CZ)** [10]: This algorithm adaptively creates balls over the joint action and context space, calculates an index for each ball based on the history of selections of that ball, and at each time step selects a matching action according to the ball with the highest index that contains the current context.

### C. Yahoo! Today Module Simulations

In YTM each instance (user) has two contexts  $(x_1, x_2) \in [0, 1]^2$ . We simulate the algorithms in Section VII-B for three different context sets in which the learning algorithms only decide based on (i) the first context  $x_1$ , (ii) the second context  $x_2$ , and (iii) both contexts  $(x_1, x_2)$  of the users. The  $m_T$  parameter of DISCOM for these simulations is set to the optimal value found in Theorem 1 (for  $\gamma = 1$ ) which is  $\lceil T^{1/4} \rceil$  for simulations with a single context and  $\lceil T^{1/5} \rceil$  for simulations with both contexts. DISCOM is run for numerous  $z$  values ranging from 1/4 to 1/2. Table IV compares the

<sup>9</sup>The number of trainings and explorations required in the regret bounds are the worst-case numbers. In reality, good performance is achieved with a much smaller number of trainings and explorations.

Context	DISCOM	DISCOM (diff. serv.)	LinUCB	Hybrid- $\epsilon$	CZ
$x_1$	6.37	7.30	6.31	5.92	4.29
$x_2$	6.14	6.45	4.72	6.14	4.39
$(x_1, x_2)$	5.93	6.61	5.65	6.15	4.24

TABLE IV

COMPARISON OF THE  $\text{CTR} \times 10^2$  ACHIEVED BY CA 1 FOR DISCOM AND OTHER LEARNING ALGORITHMS FOR YTM.

$z$	1/4	1/3	1/2
$\text{CTR} \times 10^2$	5.13	5.29	6.14
$\text{CTR} \times 10^2$ in exploitations	5.14	5.34	6.45
Exploit %	98.9	97.7	90.2
Explore %	0.5	0.6	1.9
Train %	0.7	1.7	7.9

TABLE V

THE CTR, TRAINING, EXPLORATION AND EXPLOITATION PERCENTAGES OF CA 1 USING DISCOM WITH CONTEXT  $x_2$  FOR YTM.

performance of DISCOM, LinUCB, Hybrid- $\epsilon$  and CZ. All of the algorithms are evaluated at the parameter values in which they perform the best. As seen from the table the CTR for DISCOM with differential services is 16%, 5% and 7% higher than the best of LinUCB, Hybrid- $\epsilon$  and CZ for contexts  $x_1$ ,  $x_2$  and  $(x_1, x_2)$ , respectively.

Table V compares the performance of DISCOM, the percentage of training, exploration and exploitation phases for different control functions (different  $z$  parameters) for simulations with context  $x_2$ . As expected, the percentage of trainings and explorations increase with the control function. As  $z$  increases matching actions are explored with a higher accuracy, and hence the average exploitation reward (CTR) increases.

#### D. Music Dataset Simulations

Table VI compares the performance of DISCOM, LinUCB, Hybrid- $\epsilon$  and CZ for the music dataset. The parameter values used for DISCOM for the result in Table VI are  $z = 1/8$  and  $m_T = 4$ . From the results it is observed that DISCOM achieves 10% improvement over LinUCB, 5% improvement over Hybrid- $\epsilon$ , and 28% improvement over CZ in terms of the average number of likes achieved for the users of CA 1. Moreover, the average number of likes received by DISCOM for the high type users (differential services) is even higher, which is 13%, 8% and 32% higher than LinUCB, HE and CZ, respectively.

#### E. Yahoo! Today Module with Drift Simulations

Table VII compares the performance of DISCOM-W with half window length ( $\tau_h = 2500$ ) and  $m_T = 10$ , DISCOM (with  $m_T$  set equal to  $\lceil T^{1/4} \rceil$  simulations with a single context dimension and  $\lceil T^{1/5} \rceil$  for the simulation with two context dimensions), LinUCB, Hybrid- $\epsilon$  and CZ. For the results in the table, the  $z$  parameter value of DISCOM and DISCOM-W are set to the  $z$  value in which they achieve the highest number of clicks. Similarly, LinUCB, Hybrid- $\epsilon$  and CZ are

Algorithm	DISCOM	DISCOM (diff. serv.)	LinUCB	Hybrid- $\epsilon$	CZ
Avg. num. of likes	0.717	0.736	0.652	0.683	0.559

TABLE VI

COMPARISON AMONG DISCOM AND OTHER LEARNING ALGORITHMS FOR MD.

Contexts used	DISCOM-W (diff. serv.)	DISCOM (diff. serv.)	LinUCB	Hybrid- $\epsilon$	CZ
$x_1$	6.3	5.5	5.1	3.0	2.4
$x_2$	5.1	4.2	4.2	4.6	2.4
$(x_1, x_2)$	6.9	3.8	4.6	4.1	2.3

TABLE VII

THE  $\text{CTR} \times 10^2$  OF DISCOM-W AND DISCOM FOR DIFFERENTIAL SERVICES, AND THE CTR OF OTHER LEARNING ALGORITHMS FOR YTMD.

also evaluated at their best parameter values. The results show the performance of DISCOM and DISCOM-W for differential services. DISCOM-W performs the best in this dataset in terms of the average number of clicks, with about 23%, 11.3% and 51.6% improvement over the best of LinUCB, Hybrid- $\epsilon$  and CZ, for types of contexts  $x_1$ ,  $x_2$ , and  $(x_1, x_2)$ , respectively.

## VIII. CONCLUSION

In this paper we considered novel online learning algorithms for content matching by a distributed set of CAs. We have characterized the relation between the user and content characteristics in terms of a relevance score, and proposed online learning algorithms that learn to match each user with the content with the highest relevance score. When the user and content characteristics are static, the best matching between content and each type of user can be learned perfectly, i.e., the average regret due to suboptimal matching goes to zero. When the user and content characteristics are dynamic, depending on the rate of the change, an approximately optimal matching between content and each user type can be learned. In addition to our theoretical results, we have validated the concept of distributed content matching on real-world datasets. An interesting future research direction is to investigate the interaction between different CAs when they compete for the same pool of users. Should a CA send a content that has a high chance of being liked by another CA's user to increase its immediate reward, or should it send a content that has a high chance of being disliked by the other CA's user to divert that user from using that CA and switch to it instead.

### APPENDIX A

#### A BOUND ON DIVERGENT SERIES

For  $p > 0$ ,  $p \neq 1$ ,

$$\sum_{t=1}^T t^{-p} \leq 1 + (T^{1-p} - 1)/(1 - p).$$

*Proof:* See [32]. ■

### APPENDIX B

#### PROOF OF THEOREM 1

##### A. Necessary Definitions and Notations

Let  $\beta_a := \sum_{t=1}^{\infty} 1/t^a$ , and let  $\log(\cdot)$  denote logarithm in base  $e$ . For each hypercube  $p \in \mathcal{P}_T$  let  $\bar{\pi}_{c,p} := \sup_{x \in p} \pi_c(x)$ ,  $\underline{\pi}_{c,p} := \inf_{x \in p} \pi_c(x)$ , for  $c \in \mathcal{C}$ , and  $\bar{\mu}_{k,p}^i := \sup_{x \in p} \mu_k^i(x)$ ,  $\underline{\mu}_{k,p}^i := \inf_{x \in p} \mu_k^i(x)$ , for  $k \in \mathcal{K}_i$ . Let  $x_p^*$  be the context at the center (center of symmetry) of the hypercube  $p$ . We define the optimal matching action of CA  $i$  for hypercube  $p$  as  $k_i^*(p) := \arg \max_{k \in \mathcal{K}_i} \mu_k^i(x_p^*)$ . When the hypercube  $p$

is clear from the context, we will simply denote the optimal matching action for hypercube  $p$  with  $k_i^*$ . Let

$$\mathcal{L}_p^i(t) := \left\{ k \in \mathcal{K}_i : \underline{\mu}_{k_i^*(p),p}^i - \bar{\mu}_{k,p}^i > (4Ld^{\gamma/2} + 6)t^{-z/2} \right\}$$

be the set of suboptimal matching actions of CA  $i$  at time  $t$  in hypercube  $p$ . Also related to this let

$$\mathcal{C}_p^j(t) := \left\{ c \in \mathcal{C}_j : \underline{\pi}_{c_j^*(p),p} - \bar{\pi}_{c,p} > (4Ld^{\gamma/2} + 6)t^{-z/2} \right\}$$

be the set of suboptimal contents of CA  $j$  at time  $t$  in hypercube  $p$ , where  $c_j^*(p) = \arg \max_{c \in \mathcal{C}_j} \pi_c(x_p^*)$ . Also when the hypercube  $p$  is clear from the context we will just use  $c_j^*$ . The contents in  $\mathcal{C}_p^j(t)$  are the ones that CA  $j$  should not select when called by another CA. The regret given in (3) can be written as a sum of three components:

$$R_i(T) = \mathbb{E}[R_i^e(T)] + \mathbb{E}[R_i^s(T)] + \mathbb{E}[R_i^n(T)]$$

where  $R_i^e(T)$  is the regret due to trainings and explorations by time  $T$ ,  $R_i^s(T)$  is the regret due to suboptimal matching action selections in exploitations by time  $T$  and  $R_i^n(T)$  is the regret due to near optimal matching action selections in exploitations by time  $T$ , which are all random variables.

### B. Bounding the Regret in Training, Exploration and Exploitation phases.

In the following lemmas we will bound each of these terms separately. The following lemma bounds  $\mathbb{E}[R_i^e(T)]$ .

*Lemma 1:* Consider all CAs running DISCOM with parameters  $H_1(t) = t^z \log t$ ,  $H_2(t) = C_{\max} t^z \log t$ ,  $H_3(t) = t^z \log t$  and  $m_T = \lceil T^\kappa \rceil$ , where  $0 < z < 1$  and  $0 < \kappa < 1/d$ . Then, we have

$$\mathbb{E}[R_i^e(T)] \leq 2^{d+1}(|\mathcal{C}_i| + (M-1)(C_{\max} + 1))T^{z+\kappa d} \log T + 2^{d+1}(|\mathcal{C}_i| + 2(M-1))T^{\kappa d}.$$

*Proof:* Since time slot  $t$  is a training or an exploration slot for CA  $i$  if and only if

$$\mathcal{M}_{i,p_i(t)}^{\text{ut}}(t) \cup \mathcal{M}_{i,p_i(t)}^{\text{ue}}(t) \cup \mathcal{C}_{i,p_i(t)}^{\text{ue}}(t) \neq \emptyset$$

up to time  $T$ , there can be at most  $\lceil T^z \log T \rceil$  exploration slots in which a content  $c \in \mathcal{C}_i$  is matched with the user of CA  $i$ ,  $\lceil C_{\max} T^z \log T \rceil$  training slots in which CA  $i$  selects CA  $j \in \mathcal{M}_{-i}$ ,  $\lceil T^z \log T \rceil$  exploration slots in which CA  $i$  selects CA  $j \in \mathcal{M}_{-i}$ . Result follows from summing these terms and the fact that  $(m_T)^d \leq 2^d T^{\kappa d}$  for any  $T \geq 1$ . The additional factor of 2 comes from the fact that the realized regret at any time slot can be at most 2. ■

For any  $k \in \mathcal{K}_i$  and  $p \in \mathcal{P}_T$ , the sample mean  $\bar{r}_{k,p}^i(t)$  of the relevance score of matching action  $k$  represents a random variable which is the average of the independent samples in set  $\mathcal{E}_{k,p}^i(t)$ . Since these samples are not identically distributed, in order to facilitate our analysis of the regret, we generate two different artificial i.i.d. processes to bound the probabilities related to  $\hat{\mu}_{k,p}^i(t) = \bar{r}_{k,p}^i(t) - d_k^i$ ,  $k \in \mathcal{K}_i$ . The first one is the *best* process for CA  $i$  in which the net reward of the matching action  $k$  for a user with context in  $p$  is sampled from an i.i.d. Bernoulli process with mean  $\bar{\mu}_{k,p}^i$ , the other one is the *worst* process for CA  $i$  in which this net reward is sampled from

an i.i.d. Bernoulli process with mean  $\underline{\mu}_{k,p}^i$ . Let  $\hat{\mu}_{k,p}^{\text{b},i}(z)$  denote the sample mean of the  $z$  samples from the best process and  $\hat{\mu}_{k,p}^{\text{w},i}(z)$  denote the sample mean of the  $z$  samples from the worst process for CA  $i$ . We will bound the terms  $\mathbb{E}[R_i^n(T)]$  and  $\mathbb{E}[R_i^s(T)]$  by using these artificial processes along with the similarity information given in Assumption 1.

Let  $\Xi_{j,p}^i(t)$  be the event that a suboptimal content  $c \in \mathcal{C}_j$  is selected by CA  $j \in \mathcal{M}_{-i}$ , when it is called by CA  $i$  for a context in set  $p$  for the  $t$ th time in the exploitation phases of CA  $i$ . Let  $X_{j,p}^i(t)$  denote the random variable which is the number of times CA  $j$  selects a suboptimal content when called by CA  $i$  in exploitation slots of CA  $i$  when the context is in set  $p \in \mathcal{P}_T$  by time  $t$ . Clearly, we have

$$X_{j,p}^i(t) = \sum_{t'=1}^{|\mathcal{E}_{j,p}^i(t)|} \mathbb{I}(\Xi_{j,p}^i(t'))$$

where  $\mathbb{I}(\cdot)$  is the indicator function which is equal to 1 if the event inside is true and 0 otherwise. The following lemma bounds  $\mathbb{E}[R_i^s(T)]$ .

*Lemma 2:* Consider all CAs running DISCOM with parameters  $H_1(t) = t^z \log t$ ,  $H_2(t) = C_{\max} t^z \log t$ ,  $H_3(t) = t^z \log t$  and  $m_T = \lceil T^\kappa \rceil$ , where  $0 < z < 1$  and  $\kappa = z/(2\gamma)$ . Then, we have

$$\mathbb{E}[R_i^s(T)] \leq 4(|\mathcal{C}_i| + M)\beta_2 + 4(|\mathcal{C}_i| + M)MC_{\max}\beta_2 \frac{T^{1-z/2}}{1-z/2}.$$

*Proof:* Consider time  $t$ . For simplicity of notation let  $p = p_i(t)$ . Let

$$\mathcal{W}^i(t) := \{\mathcal{M}_{i,p_i(t)}^{\text{ut}}(t) \cup \mathcal{M}_{i,p_i(t)}^{\text{ue}}(t) \cup \mathcal{C}_{i,p_i(t)}^{\text{ue}}(t) = \emptyset\}$$

be the event that CA  $i$  exploits at time  $t$ .

First, we will bound the probability that CA  $i$  selects a suboptimal matching action in an exploitation slot. Then, using this we will bound the expected number of times a suboptimal matching action is selected by CA  $i$  in exploitation slots. Note that every time a suboptimal matching action is selected by CA  $i$ , since  $\mu_k^i(x) = \pi_k^i(x) - d_k^i \in [-1, 1]$  for all  $k \in \mathcal{K}_i$ , the realized (hence expected) loss is bounded above by 2. Therefore 2 times the expected number of times a suboptimal matching action is chosen in an exploitation slot bounds the regret due to suboptimal matching actions in exploitation slots. Let  $\mathcal{V}_k^i(t)$  be the event that matching action  $k \in \mathcal{K}_i$  is chosen at time  $t$  by CA  $i$ . We have

$$R_i^s(T) \leq 2 \sum_{t=1}^T \sum_{k \in \mathcal{L}_{p_i(t)}^i} \mathbb{I}(\mathcal{V}_k^i(t), \mathcal{W}^i(t)).$$

Taking the expectation

$$\mathbb{E}[R_i^s(T)] \leq 2 \sum_{t=1}^T \sum_{k \in \mathcal{L}_{p_i(t)}^i} \mathbb{P}(\mathcal{V}_k^i(t), \mathcal{W}^i(t)). \quad (10)$$

Let  $\mathcal{B}_{j,p}^i(t)$  be the event that at most  $t^\phi$  samples in  $\mathcal{E}_{j,p}^i(t)$  are collected from suboptimal content of CA  $j$ . Let  $\mathcal{B}^i(t) := \bigcap_{j \in \mathcal{M}_{-i}} \mathcal{B}_{j,p_i(t)}^i(t)$ . For a set  $\mathcal{A}$ , let  $\mathcal{A}^c$  denote the complement

of that set. For any  $k \in \mathcal{K}_i$ , we have

$$\begin{aligned}
& \mathbb{P}(\mathcal{V}_{k,p}^i(t), \mathcal{W}^i(t)) \\
& \leq \mathbb{P}(\hat{\mu}_{k,p}^i(t) \geq \bar{\mu}_{k,p}^i + H_t, \mathcal{W}^i(t), \mathcal{B}^i(t)) \\
& + \mathbb{P}(\hat{\mu}_{k_i^*,p}^i(t) \leq \underline{\mu}_{k_i^*,p}^i - H_t, \mathcal{W}^i(t), \mathcal{B}^i(t)) + \mathbb{P}(\mathcal{B}^i(t)^c) \\
& + \mathbb{P}(\hat{\mu}_{k,p}^i(t) \geq \hat{\mu}_{k_i^*,p}^i(t), \hat{\mu}_{k,p}^i(t) < \bar{\mu}_{k,p}^i + H_t, \\
& \hat{\mu}_{k_i^*,p}^i(t) > \underline{\mu}_{k_i^*,p}^i - H_t, \mathcal{W}^i(t), \mathcal{B}^i(t)) \quad (11)
\end{aligned}$$

for some  $H_t > 0$ . We have for any suboptimal matching action  $k \in \mathcal{L}_p^i(t)$ ,

$$\begin{aligned}
& \mathbb{P}(\hat{\mu}_{k,p}^i(t) \geq \hat{\mu}_{k_i^*,p}^i(t), \hat{\mu}_{k,p}^i(t) < \bar{\mu}_{k,p}^i + H_t, \\
& \hat{\mu}_{k_i^*,p}^i(t) > \underline{\mu}_{k_i^*,p}^i - H_t, \mathcal{W}^i(t), \mathcal{B}_{k,p}^i(t)) \\
& \leq \mathbb{P}(\hat{\mu}_{k,p}^{b,i}(|\mathcal{E}_{k,p}^i(t)|) \geq \hat{\mu}_{k_i^*,p}^{w,i}(|\mathcal{E}_{k_i^*,p}^i(t)|) - t^{\phi-z}, \\
& \hat{\mu}_{k,p}^{b,i}(|\mathcal{E}_{k,p}^i(t)|) < \bar{\mu}_{k,p}^i + L(\sqrt{d}/m_T)^\gamma + H_t + t^{\phi-z}, \\
& \hat{\mu}_{k_i^*,p}^{w,i}(|\mathcal{E}_{k_i^*,p}^i(t)|) > \underline{\mu}_{k_i^*,p}^i - L(\sqrt{d}/m_T)^\gamma - H_t, \mathcal{W}^i(t)).
\end{aligned}$$

For  $k \in \mathcal{L}_p^i(t)$ , when

$$2L(\sqrt{d}/m_T)^\gamma + 2H_t + 2t^{\phi-z} \leq (4Ld^{\gamma/2} + 6)t^{-z/2} \quad (12)$$

the three inequalities given below

$$\begin{aligned}
& \underline{\mu}_{k_i^*,p}^i - \bar{\mu}_{k,p}^i > (4Ld^{\gamma/2} + 6)t^{-z/2} \\
& \hat{\mu}_{k,p}^{b,i}(|\mathcal{E}_{k,p}^i(t)|) < \bar{\mu}_{k,p}^i + L(\sqrt{d}/m_T)^\gamma + H_t + t^{\phi-z} \\
& \hat{\mu}_{k_i^*,p}^{w,i}(|\mathcal{E}_{k_i^*,p}^i(t)|) > \underline{\mu}_{k_i^*,p}^i - L(\sqrt{d}/m_T)^\gamma - H_t
\end{aligned}$$

together imply that  $\hat{\mu}_{k,p}^{b,i}(|\mathcal{E}_{k,p}^i(t)|) < \hat{\mu}_{k_i^*,p}^{w,i}(|\mathcal{E}_{k_i^*,p}^i(t)|) - t^{\phi-z}$ , which implies that

$$\begin{aligned}
& \mathbb{P}(\hat{\mu}_{k,p}^i(t) \geq \hat{\mu}_{k_i^*,p}^i(t), \hat{\mu}_{k,p}^i(t) < \bar{\mu}_{k,p}^i + H_t, \\
& \hat{\mu}_{k_i^*,p}^i(t) > \underline{\mu}_{k_i^*,p}^i - H_t, \mathcal{W}^i(t), \mathcal{B}_{k,p}^i(t)) = 0. \quad (13)
\end{aligned}$$

Let  $H_t = 2t^{\phi-z} + Ld^{\gamma/2}m_T^{-\gamma}$ . A sufficient condition that implies (12) is

$$4Ld^{\gamma/2}t^{-\kappa\gamma} + 6t^{\phi-z} \leq (4Ld^{\gamma/2} + 6)t^{-z/2} \quad (14)$$

which holds for all  $t \geq 1$  when  $\phi = z/2$  and  $\kappa\gamma \geq z/2$ . Using a Chernoff-Hoeffding bound, for any  $k \in \mathcal{L}_{p_i(t)}^i(t)$ , since on the event  $\mathcal{W}^i(t)$ ,  $|\mathcal{E}_{k,p_i(t)}^i(t)| \geq t^z \log t$ , we have

$$\mathbb{P}(\hat{\mu}_{k,p}^i(t) \geq \bar{\mu}_{k,p}^i + H_t, \mathcal{W}^i(t), \mathcal{B}^i(t)) \leq t^{-2} \quad (15)$$

and

$$\mathbb{P}(\hat{\mu}_{k_i^*,p}^i(t) \leq \underline{\mu}_{k_i^*,p}^i - H_t, \mathcal{W}^i(t), \mathcal{B}^i(t)) \leq t^{-2}. \quad (16)$$

Since  $\{\mathcal{B}_{j,p}^i(t)^c, \mathcal{W}^i(t)\} = \{X_{j,p}^i(t) \geq t^\phi\}$ , by applying the Markov inequality, we have

$$\mathbb{P}(\mathcal{B}_{j,p}^i(t)^c, \mathcal{W}^i(t)) \leq \mathbb{E}[X_{j,p}^i(t)]t^{-\phi}.$$

Since

$$X_{j,p}^i(t) = \sum_{t'=1}^{|\mathcal{E}_{j,p}^i(t)|} \mathbb{I}(\Xi_{j,p}^i(t'))$$

and

$$\begin{aligned}
& \mathbb{P}(\Xi_{j,p}^i(t)) \\
& \leq \sum_{m \in \mathcal{C}_p^j(t)} \mathbb{P}(\bar{r}_{m,p}^j(t) \geq \bar{r}_{c_j^*,p}^j(t)) \\
& \leq \sum_{m \in \mathcal{C}_p^j(t)} (\mathbb{P}(\bar{r}_{m,p}^j(t) \geq \bar{\pi}_{m,p} + H_t, \mathcal{W}^i(t)) \\
& + \mathbb{P}(\bar{r}_{c_j^*,p}^j(t) \leq \underline{\pi}_{c_j^*,p} - H_t, \mathcal{W}^i(t)) + \mathbb{P}(\bar{r}_{m,p}^j(t) \geq \bar{r}_{c_j^*,p}^j(t), \\
& \bar{r}_{m,p}^j(t) < \bar{\pi}_{m,p} + H_t, \bar{r}_{c_j^*,p}^j(t) > \underline{\pi}_{c_j^*,p} - H_t, \mathcal{W}^i(t))).
\end{aligned}$$

When (14) holds, the last probability in the sum above is equal to zero while the first two probabilities are upper bounded by  $e^{-2(H_t)^2 t^z \log t}$ . Thus, we have

$$\mathbb{P}(\Xi_{j,p}^i(t)) \leq \sum_{m \in \mathcal{C}_p^j(t)} 2e^{-2(H_t)^2 t^z \log t} \leq 2|\mathcal{C}_j|t^{-2}.$$

This implies that

$$\mathbb{E}[X_{j,p}^i(t)] \leq \sum_{t'=1}^{\infty} \mathbb{P}(\Xi_{j,p}^i(t')) \leq 2|\mathcal{C}_j| \sum_{t'=1}^{\infty} (t')^{-2}.$$

Therefore, by the Markov inequality and union bound we get

$$\begin{aligned}
& \mathbb{P}(\mathcal{B}_{j,p_i(t)}^i(t)^c, \mathcal{W}^i(t)) = \mathbb{P}(X_{j,p_i(t)}^i(t) \geq t^\phi) \\
& \leq 2|\mathcal{C}_j|\beta_2 t^{-z/2}
\end{aligned}$$

and

$$\mathbb{P}(\mathcal{B}^i(t)^c, \mathcal{W}^i(t)) \leq 2MC_{\max}\beta_2 t^{-z/2}. \quad (17)$$

Then, using (13), (15), (16) and (17), we have

$$\mathbb{P}(\mathcal{V}_k^i(t), \mathcal{W}^i(t)) \leq 2t^{-2} + 2MC_{\max}\beta_2 t^{-z/2},$$

for any  $k \in \mathcal{L}_{p_i(t)}^i(t)$ , and By (10), and by the result of Appendix A, we get the stated bound for  $\mathbb{E}[R_i^s(T)]$ . ■

The next lemma bounds  $\mathbb{E}[R_i^n(T)]$ .

*Lemma 3:* Consider all CAs running DISCOM with parameters  $H_1(t) = t^z \log t$ ,  $H_2(t) = C_{\max}t^z \log t$ ,  $H_3(t) = t^z \log t$  and  $m_T = \lceil T^\kappa \rceil$ , where  $0 < z < 1$  and  $\kappa = z/(2\gamma)$ . Then, we have

$$\mathbb{E}[R_i^n(T)] \leq \frac{(14Ld^{\gamma/2} + 12)}{1 - z/2} T^{1-z/2} + 4C_{\max}\beta_2.$$

*Proof:* At any time  $t$ , for any  $k \in \mathcal{K}_i - \mathcal{L}_p^i(t)$  and  $x \in p$ , we have

$$\mu_{k_i^*,p}^i(x) - \mu_k^i(x) \leq (7Ld^{\gamma/2} + 6)t^{-z/2}.$$

Similarly, for any  $j \in \mathcal{M}$ ,  $c \in \mathcal{C}_j - \mathcal{C}_p^j(t)$  and  $x \in p$ , we have

$$\pi_{c_j^*,p}^i(x) - \pi_c(x) \leq (7Ld^{\gamma/2} + 6)t^{-z/2}.$$

Due to the above inequalities, if a near optimal action in  $\mathcal{C}_i \cap (\mathcal{K}_i - \mathcal{L}_p^i(t))$  is chosen by CA  $i$  at time  $t$ , the contribution to the regret is at most  $(7Ld^{\gamma/2} + 6)t^{-z/2}$ . If a near optimal



CA  $j \in \mathcal{M}_{-i} \cap (\mathcal{K}_i - \mathcal{L}_p^i(t))$  is called by CA  $i$  at time  $t$ , and if CA  $j$  selects one of its near optimal contents in  $\mathcal{C}_j - \mathcal{C}_p^j(t)$ , then the contribution to the regret is at most  $2(7Ld^{\gamma/2} + 6)t^{-z/2}$ . Moreover since we are in an exploitation step, the near-optimal CA  $j$  that is chosen can choose one of its suboptimal contents in  $\mathcal{C}_p^j(t)$  with probability at most  $2C_{\max}t^{-2}$ , which will result in an expected regret of at most  $4C_{\max}t^{-2}$ .

Therefore, the total regret due to near optimal choices of CA  $i$  by time  $T$  is upper bounded by

$$\begin{aligned} & (14Ld^{\gamma/2} + 12) \sum_{t=1}^T t^{-z/2} + 4C_{\max} \sum_{t=1}^T t^{-2} \\ & \leq \frac{(14Ld^{\gamma/2} + 12)}{1 - z/2} T^{1-z/2} + 4C_{\max}\beta_2. \end{aligned}$$

by using the result in Appendix A.  $\blacksquare$

Next, we give the proof of Theorem 1 by combining the results of the above lemmas.

### C. Proof of Theorem 1

The highest orders of regret that come from Lemmas 1, 2 and 3 are  $\tilde{O}(T^{\kappa d + z})$ ,  $O(T^{1-z/2})$ ,  $O(T^{1-z/2})$ . We need to optimize them with respect to the constraint Regret is minimized when  $\kappa d + z = 1 - z/2$ , which is attained by  $z = 2\gamma/(3\gamma + d)$ . Result follows from summing the bounds in Lemmas 1, 2 and 3.

#### APPENDIX C PROOF OF COROLLARY 1

From the proof of Lemma 2, for  $z = 2\gamma/(3\gamma + d)$ , we have

$$\mathbb{P}(\mathcal{V}_k^i(t), \mathcal{W}^i(t)) \leq 2t^{-2} + 2MC_{\max}\beta_2 t^{-\gamma/(3\gamma+d)}$$

for  $k \in \mathcal{L}_{p_i(t)}^i(t)$ . This implies that

$$\begin{aligned} & \mathbb{P}(a_i(t) \in \mathcal{L}_{p_i(t)}^i(t), \mathcal{W}^i(t)) \\ & \leq \sum_{k \in \mathcal{L}_{p_i(t)}^i(t)} \mathbb{P}(\mathcal{V}_k^i(t), \mathcal{W}^i(t)) \\ & \leq \frac{2|\mathcal{K}_i|}{t^2} + \frac{2|\mathcal{K}_i|MC_{\max}\beta_2}{t^{\gamma/(3\gamma+d)}}. \end{aligned}$$

The difference between the expected reward of an action within a hypercube from its expected reward at the center of the hypercube is at most  $Ld^{\gamma/2}/(m_T)^\gamma$ . Since  $m_T = \lceil T^{1/(3\gamma+d)} \rceil$ ,  $a_i(t) \in \mathcal{K}_i - \mathcal{L}_{p_i(t)}^i(t)$  implies that

$$\mu_{a_i(t)}^i(x_i(t)) \geq \mu_{k_i^*(x_i(t))}^i(x_i(t)) - (6Ld^{\gamma/2} + 6)T^{-\gamma/(3\gamma+d)}.$$

#### APPENDIX D PROOF OF THEOREM 2

In order for time  $t$  to be an exploitation slot for CA  $i$  it is required that  $\mathcal{M}_{i,p_i(t)}^{\text{ut}}(t) \cup \mathcal{M}_{i,p_i(t)}^{\text{uc}}(t) \cup \mathcal{C}_{i,p_i(t)}^{\text{uc}}(t) = \emptyset$ . Since the counters of DISCOM are updated only when feedback is received, and since the control functions are the same as the ones that are used in the setting where feedback is always available, the regret due to suboptimal and near optimal matching actions by time  $t$  with missing feedback will not be any greater than the regret due to suboptimal

and near optimal matching actions for the case when the users always provide feedback. Therefore, the bounds given in Lemmas 2 and 3 will also hold for the case with missing feedback. Only the regret due to trainings and explorations increases, since more trainings and explorations are needed before the counters exceed the values of the control functions such that the relevance score estimates are accurate enough to exploit. Consider any  $p \in \mathcal{P}_T$ . From the definition of DISCOM, the number of exploration slots in which content  $c \in \mathcal{C}_i$  is matched with CA  $i$ 's user and the user's feedback is observed is at most  $\lceil T^{2\gamma/(3\gamma+d)} \rceil$ . The number of training slots in which CA  $i$  requested content from CA  $j \in \mathcal{M}_{-i}$  and received the feedback about this content from its user is at most  $\lceil C_{\max}T^{2\gamma/(3\gamma+d)} \log T \rceil$ . The number of exploration slots in which CA  $i$  selected CA  $j \in \mathcal{M}_{-i}$  is at most  $\lceil T^{2\gamma/(3\gamma+d)} \log T \rceil$ .

Let  $\tau_{\text{exp}}(T)$  be the random variable which denotes the smallest time step for which for each  $c \in \mathcal{C}_i$  there are  $\lceil T^{2\gamma/(3\gamma+d)} \rceil$  feedback observations, for each  $j \in \mathcal{M}_{-i}$  there are  $\lceil C_{\max}T^{2\gamma/(3\gamma+d)} \log T \rceil$  feedback observations for the trainings and  $\lceil T^{2\gamma/(3\gamma+d)} \log T \rceil$  feedback observations for the explorations. Then,  $\mathbb{E}[\tau_{\text{exp}}(T)]$  is the expected number of training plus exploration slots by time  $T$ . Let  $Y_{\text{exp}}(t)$  be the random variable which denotes the number of time slots in which the feedback is not provided by the users of CA  $i$  till CA  $i$  received  $t$  feedbacks from its users. Let  $A_i(T) = Z_i T^{2\gamma/(3\gamma+d)} \log T + (|\mathcal{C}_i| + 2(M-1))$ . We have

$$\mathbb{E}[\tau_{\text{exp}}(T)] = \mathbb{E}[Y_{\text{exp}}(A_i(T))] + A_i(T).$$

$Y_{\text{exp}}(A_i(T))$  is a negative binomial random variable with probability of observing no feedback at any time  $t$  equals to  $1 - p_r$ . Therefore,

$$\mathbb{E}[Y_{\text{exp}}(A_i(T))] = (1 - p_r)A_i(T)/p_r.$$

Using this, we get

$$\mathbb{E}[\tau_{\text{exp}}(T)] = A_i(T)/p_r.$$

The regret bound follows from substituting this into the proof of Theorem 1.

#### APPENDIX E PROOF OF THEOREM 3

The basic idea is to choose  $\tau_h$  in a way that the regret due to variation of relevance scores over time and the regret due to variation of estimated relevance scores due to the limited number of observations during each round is balanced. Majority of the steps of this proof is similar to the proof of Theorem 1 hence some of the steps are omitted.

Consider a round  $\eta$  of length  $2\tau_h$ . Denote the set of time slots in round  $\eta$  by  $[\eta]$ . For any  $c \in \mathcal{C}$  let

$$\bar{\pi}_{c,p,\eta} := \sup_{x \in p, t \in [\eta]} \pi_{c,t}(x),$$

$$\underline{\pi}_{c,p,\eta} := \inf_{x \in p, t \in [\eta]} \pi_{c,t}(x).$$

For any  $k \in \mathcal{K}_i$  let

$$\bar{\mu}_{k,p,\eta}^i := \sup_{x \in p, t \in [\eta]} \mu_{k,t}^i(x),$$

$$\underline{\mu}_{k,p,\eta}^i := \inf_{x \in p, t \in [\eta]} \mu_{k,t}^i(x),$$

where  $\mu_{k,t}^i(x)$  is defined as the time-varying version of  $\mu_k^i(x)$  given in (1) under Assumption 2. For CA  $i$ , the set of suboptimal matching actions is given as

$$\begin{aligned} \mathcal{L}_{p,\eta}^i(t) &:= \left\{ k \in \mathcal{K}_i : \underline{\mu}_{k_i^*(p),p,\eta}^i - \bar{\mu}_{k,p,\eta}^i \right. \\ &> (4Ld^{\gamma/2} + 6)(t \bmod \tau_h + 1)^{-z/2} + \frac{4\tau_h}{T_s} \left. \right\}, \end{aligned}$$

where  $k_i^*(p)$  is the matching action with the highest net reward for the context at the center of  $p$  at the time slot in the middle of round  $\eta$ .

Consider the regret due to explorations and trainings for DISCOM $_{\eta}$  incurred over times when it is in the active sub-phase (over  $\tau_h$  time slots). Similar to the proof of Lemma 1 it can be shown that the regret due to trainings and explorations is

$$\mathbb{E}[R_i^e(\tau_h)] = \tilde{O}(\tau_h^{z+\kappa d}).$$

Similar to the proof of Lemma 2, it can be shown that the regret due to suboptimal matching action selections is

$$\mathbb{E}[R_i^s(\tau_h)] = O\left(\tau_h^{1-z/2}\right)$$

when  $\kappa = z/(2\gamma)$ . Since the definition of a sub-optimal matching action is different for dynamic user and content characteristics, the regret due to near optimal matching actions in  $\mathcal{K}_i - \mathcal{L}_{p,\eta}^i(t)$  is different from Lemma 3. At time  $t$  which is in round  $\eta$ , since a near optimal matching action's contribution to the one-step regret is at most

$$(8Ld^{\gamma/2} + 12)(t \bmod \tau_h + 1)^{-z/2} + 4\tau_h/T_s$$

summing over all time slots in a round  $\eta$ , we have

$$\mathbb{E}[R_i^n(\tau_h)] = O\left(\tau_h^{1-z/2}\right) + O\left(\frac{\tau_h^2}{T_s}\right).$$

Clearly we have  $\mathbb{E}[R_i^s(\tau_h)] \leq \mathbb{E}[R_i^e(\tau_h)]$ . Let  $\tau_h = \lfloor T_s^\phi \rfloor$  for some  $\phi > 0$ . Then we have

$$\frac{\mathbb{E}[R_i^e(\tau_h)]}{\tau_h} = \tilde{O}\left(T_s^{\phi z + \phi \kappa d - \phi}\right),$$

and

$$\frac{\mathbb{E}[R_i^n(\tau_h)]}{\tau_h} = O\left(T_s^{-\phi z/2}\right) + O\left(T_s^{\phi-1}\right).$$

The sum  $(\mathbb{E}[R_i^e(\tau_h)] + \mathbb{E}[R_i^s(\tau_h)] + \mathbb{E}[R_i^n(\tau_h)])/\tau_h$  is minimized by setting  $z = 2\gamma/(3\gamma + d)$  and  $\phi = 1/(1 + z/2)$ . Hence,  $\tau_h = \lfloor T_s^{\frac{\gamma}{4\gamma+d}} \rfloor$  the order of the time averaged regret is equal to  $\tilde{O}\left(T_s^{\frac{\gamma}{4\gamma+d}}\right)$ .

## REFERENCES

[1] C. Tekin and M. van der Schaar, "Contextual online learning for multimedia content aggregation," to appear in *IEEE Trans. Multimedia*, 2015.  
[2] S. Ren and M. van der Schaar, "Pricing and investment for online TV content platforms," *IEEE Trans. Multimedia*, vol. 14, no. 6, pp. 1566–1578, 2012.

[3] S. Song, H. Moustafa, and H. Afifi, "Advanced IPTV services personalization through context-aware content recommendation," *IEEE Trans. Multimedia*, vol. 14, no. 6, pp. 1528–1537, 2012.  
[4] C. Xu, J. Wang, H. Lu, and Y. Zhang, "A novel framework for semantic annotation and personalized retrieval of sports video," *IEEE Trans. Multimedia*, vol. 10, no. 3, pp. 421–436, 2008.  
[5] L. Li, W. Chu, J. Langford, and R. E. Schapire, "A contextual-bandit approach to personalized news article recommendation," in *Proc. 19th Int. Conf. World Wide Web*, 2010, pp. 661–670.  
[6] M. Saxena, U. Sharan, and S. Fahmy, "Analyzing video services in web 2.0: a global perspective," in *Proc. 18th Int. Workshop Network and Operating Systems Support for Digital Audio and Video*, 2008, pp. 39–44.  
[7] R. Mohan, J. R. Smith, and C.-S. Li, "Adapting multimedia internet content for universal access," *IEEE Trans. Multimedia*, vol. 1, no. 1, pp. 104–114, 1999.  
[8] M. Schranz, S. Dustdar, and C. Platzer, "Building an integrated pan-European news distribution network," in *Collaborative Networks and Their Breeding Environments*. Springer US, 2005, pp. 587–596.  
[9] A. Boutet, K. Kloudas, and A.-M. Kermarrec, "FStream: a decentralized and social music streamer," in *Proc. Int. Conf. Networked Systems (NETYS)*, pp. 253–257.  
[10] A. Slivkins, "Contextual bandits with similarity information," in *Proc. 24th Annual Conf. Learning Theory (COLT)*, vol. 19, June 2011, pp. 679–702.  
[11] T. Lu, D. Pál, and M. Pál, "Contextual multi-armed bandits," in *Proc. 13th Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, vol. 9, May 2010, pp. 485–492.  
[12] Y. Deshpande and A. Montanari, "Linear bandits in high dimension and recommendation systems," in *Proc. 50th Annual Allerton Conf. Communication, Control and Computing*, 2012, pp. 1750–1754.  
[13] P. Kohli, M. Salek, and G. Stoddard, "A fast bandit algorithm for recommendations to users with heterogeneous tastes," in *Proc. 27th Conf. Artificial Intelligence (AAAI)*, July 2013, pp. 1135–1141.  
[14] N. Sahoo, P. V. Singh, and T. Mukhopadhyay, "A hidden Markov model for collaborative filtering," *MIS Quarterly*, vol. 36, no. 4, pp. 1329–1356, 2012.  
[15] G. Linden, B. Smith, and J. York, "Amazon.com recommendations: Item-to-item collaborative filtering," *Internet Comput.*, vol. 7, no. 1, pp. 76–80, 2003.  
[16] K. Miyahara and M. J. Pazzani, "Collaborative filtering with the simple Bayesian classifier," in *PRICAI 2000 Topics in Artificial Intelligence*. Springer, 2000, pp. 679–689.  
[17] S. D. Roy, T. Mei, W. Zeng, and S. Li, "Empowering cross-domain internet media with real-time topic learning from social streams," in *Proc. IEEE Int. Conf. Multimedia and Expo (ICME)*, 2012, pp. 49–54.  
[18] S. Roy, T. Mei, W. Zeng, and S. Li, "Towards cross-domain learning for social video popularity prediction," *IEEE Trans. Multimedia*, vol. 15, no. 6, pp. 1255–1267, Oct 2013.  
[19] D. Bouneffouf, A. Bouzeghoub, and A. L. Gañarski, "Hybrid-ε-greedy for mobile context-aware recommender system," in *Advances in Knowledge Discovery and Data Mining*. Springer, 2012, pp. 468–479.  
[20] E. Hazan and N. Megiddo, "Online learning with prior knowledge," in *Learning Theory*. Springer, 2007, pp. 499–513.  
[21] M. Dudik, D. Hsu, S. Kale, N. Karampatziakis, J. Langford, L. Reyzin, and T. Zhang, "Efficient optimal learning for contextual bandits," *arXiv preprint arXiv:1106.2369*, 2011.  
[22] W. Chu, L. Li, L. Reyzin, and R. E. Schapire, "Contextual bandits with linear payoff functions," in *Proc. 14th Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, vol. 15, April 2011, pp. 208–214.  
[23] P. Auer, N. Cesa-Bianchi, and P. Fischer, "Finite-time analysis of the multiarmed bandit problem," *Machine Learning*, vol. 47, p. 235256, 2002.  
[24] M. Naaman, "Social multimedia: highlighting opportunities for search and mining of multimedia data in social media applications," *Multimedia Tools and Applications*, vol. 56, no. 1, pp. 9–34, 2012.  
[25] L. L. Minku, A. P. White, and X. Yao, "The impact of diversity on online ensemble learning in the presence of concept drift," *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 5, pp. 730–742, 2010.  
[26] A. Narasimhamurthy and L. I. Kuncheva, "A framework for generating data to simulate changing environments," in *Proc. 25th IASTED Int. Multi-Conf.: Artificial Intelligence and Applications*, February 2007, pp. 384–389.  
[27] M. van der Schaar, J. Xu, and W. Zame, "Efficient online exchange via fiat money," *Economic Theory*, vol. 54, no. 2, pp. 211–248, 2013.  
[28] S. Nakamoto, "Bitcoin: A peer-to-peer electronic cash system," *Consulted*, vol. 1, no. 2012, p. 28, 2008.

- [29] J. Langford and T. Zhang, "The epoch-greedy algorithm for contextual multi-armed bandits," in *Proc. 20th Int. Conf. Neural Information Processing Systems (NIPS)*, vol. 20, 2007, pp. 1096–1103.
- [30] J. Gao, W. Fan, and J. Han, "On appropriate assumptions to mine data streams: Analysis and practice," in *Proc. 7th IEEE Int. Conf. Data Mining (ICDM)*, 2007, pp. 143–152.
- [31] M. M. Masud, J. Gao, L. Khan, J. Han, and B. Thuraisingham, "Integrating novel class detection with classification for concept-drifting data streams," in *Machine Learning and Knowledge Discovery in Databases*. Springer, 2009, pp. 79–94.
- [32] E. Chlebus, "An approximate formula for a partial sum of the divergent p-series," *Applied Mathematics Letters*, vol. 22, no. 5, pp. 732–737, 2009.



**Cem Tekin** Cem Tekin is an Assistant Professor in Electrical and Electronics Engineering Department at Bilkent University, Turkey. From February 2013 to January 2015, he was a Postdoctoral Scholar at University of California, Los Angeles. He received the B.Sc. degree in electrical and electronics engineering from the Middle East Technical University, Ankara, Turkey, in 2008, the M.S.E. degree in electrical engineering: systems, M.S. degree in mathematics, Ph.D. degree in electrical engineering: systems from the University of Michigan, Ann Arbor, in 2010, 2011

and 2013, respectively. His research interests include machine learning, multi-armed bandit problems, data mining, multi-agent systems and game theory. He received the University of Michigan Electrical Engineering Departmental Fellowship in 2008, and the Fred W. Ellersick award for the best paper in MILCOM 2009.



**Mihaela van der Schaar** Mihaela van der Schaar is Chancellor Professor of Electrical Engineering at University of California, Los Angeles. Her research interests include network economics and game theory, online learning, dynamic multi-user networking and communication, multimedia processing and systems, real-time stream mining. She is an IEEE Fellow, a Distinguished Lecturer of the Communications Society for 2011-2012, the Editor in Chief of IEEE Transactions on Multimedia and a member of the Editorial Board of the IEEE Journal on Selected

Topics in Signal Processing. She received an NSF CAREER Award (2004), the Best Paper Award from IEEE Transactions on Circuits and Systems for Video Technology (2005), the Okawa Foundation Award (2006), the IBM Faculty Award (2005, 2007, 2008), the Most Cited Paper Award from EURASIP: Image Communications Journal (2006), the Gamenets Conference Best Paper Award (2011) and the 2011 IEEE Circuits and Systems Society Darlington Award Best Paper Award. She received three ISO awards for her contributions to the MPEG video compression and streaming international standardization activities, and holds 33 granted US patents.