

A Concentration of Measure Approach to Correlated Graph Matching

Farhad Shirani^{ID}, Member, IEEE, Siddharth Garg, and Elza Erkip^{ID}, Fellow, IEEE

Abstract—The graph matching problem emerges naturally in various applications such as Web privacy, image processing and computational biology. In this article, graph matching is considered under a stochastic model, where a pair of randomly generated graphs with pairwise correlated edges are to be matched such that given the labeling of the vertices in the first graph, the labels in the second graph are recovered by leveraging the correlation among their edges. The problem is considered under various settings and graph models. In the first step, the Correlated Erdős-Rényi (CER) graph model is studied, where all edge pairs whose vertices have similar labels are generated based on identical distributions and independently of other edges. A matching scheme called the *typicality matching scheme* is introduced. The scheme operates by investigating the joint typicality of the adjacency matrices of the two graphs. New results on the typicality of permutations of sequences lead to necessary and sufficient conditions for successful matching based on the parameters of the CER model. In the next step, the results are extended to graphs with community structure generated based on the Stochastic Block Model (SBM). The SBM model is a generalization of the CER model where each vertex in the graph is associated with a community label, which affects its edge statistics. The results are further extended to matching of ensembles of more than two correlated graphs. Lastly, the problem of seeded graph matching is investigated where a subset of the labels in the second graph are known prior to matching. In this scenario, in addition to obtaining necessary and sufficient conditions for successful matching, a polynomial time matching algorithm is proposed.

Index Terms—Network theory, graph theory, data privacy, information theory, graph matching, graph alignment, attributed graphs, typicality matching, permuted sequences, correlated graphs.

I. INTRODUCTION

ONLINE social networks store large quantities of personal data from their users. As a result, social network privacy

Manuscript received August 16, 2020; revised December 4, 2020 and January 22, 2021; accepted January 24, 2021. Date of publication February 2, 2021; date of current version March 16, 2021. This work was supported in part by NSF under Grant CCF-1815821 and Grant CNS-1619129, and in part by the ND EPSCoR under Grant FAR0033968. This work was presented in part at IEEE International Symposium on Information Theory (ISIT), July 2018, 51st Asilomar Conference on Signals, Systems, and Computers, November 2017, and 56th Annual Allerton Conference on Communication, Control, and Computing, October 2018. (Corresponding author: Farhad Shirani.)

Farhad Shirani is with the Electrical and Computer Engineering Department, North Dakota State University, Fargo, ND 58105 USA (e-mail: f.shiranichaharsoogh@ndsu.edu).

Siddharth Garg and Elza Erkip are with the Electrical and Computer Engineering Department, New York University, New York City, NY 11201 USA (e-mail: siddharth.garg@nyu.edu; elza@nyu.edu).

Digital Object Identifier 10.1109/JSAT.2021.3056280

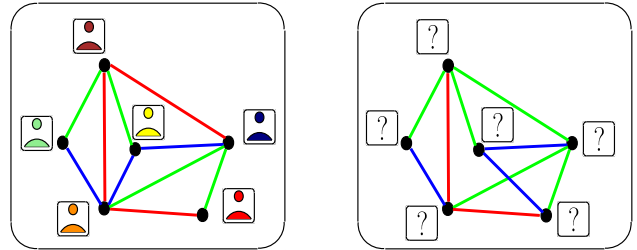


Fig. 1. An instance of the graph matching problem where the anonymized graph on the right is to be matched to the de-anonymized graph on the left.

has become an issue of significant concern. Social network data is often released to third-parties in an anonymized and obfuscated form for various purposes including targeted advertising, developing new applications, and academic research [1], [2]. However, it has been pointed out that anonymizing social network data through removing user IDs before publishing the data is far from enough to protect users' privacy [3], [4]. To elaborate, it has been shown through real-world implementation of privacy attacks that an attacker can potentially recover the user IDs by aligning the user profiles in the anonymized social network graph with the public profiles of users in other social networks on the Web. In other words, the attacker can 'match' the anonymized social network profiles of users with their public profiles in other social networks. *Graph Matching* — also known as network alignment — describes the problem of detecting node correspondence across graphs. In addition to social network deanonymization [5]–[7], the need for matching two or more graphs arises naturally in a variety of other applications of interest such as pattern recognition [8], cross-lingual knowledge alignment [9], and protein interaction network alignment [10]. The significant increase in the ability to store, share, and analyze large graphs has led to a growing need to develop *low complexity* algorithms for graph matching, and derive *theoretical guarantees* for their success, that is, to study how and when is it possible to perform fast and efficient network alignment.

In the simplest form of graph matching scenarios, an agent is given a correlated pair of randomly generated graphs: i) an 'anonymized' unlabeled graph, and ii) a 'de-anonymized' labeled graph as shown in Figure 1. The objective is to leverage the correlation among the edges of the graphs to recover the canonical labeling of the vertices in the anonymized graph. The fundamental limits of graph matching, i.e., characterizing the necessary and sufficient conditions on graph parameters for successful matching, has been considered under various

probabilistic models capturing the correlation among the graph edges. In the *Correlated Erdős-Rényi* (CER) model the edges in the two graphs are pairwise correlated and are generated independently, based on identical distributions. More precisely, in this model, edges whose vertices are labeled identically are correlated through an arbitrary joint probability distribution and are generated independently of all other edges. In its simplest form — where the edges of the two graphs are exactly equal — graph matching is called *graph isomorphism*. Tight necessary and sufficient conditions for successful matching in the graph isomorphism scenario were derived in [11], [12] and polynomial time algorithms were proposed in [13]–[15]. The problem of matching non-identical pairs of CER graphs was studied in [16]–[22] and conditions for successful matching were derived.

The CER model assumes the existence of statistical correlation among edge pairs connecting matching vertices in the two graphs, where the correlation model is based on an identical distribution among all matching edge pairs. Consequently, it does not model the community structure among the graph nodes which manifests in many applications [23], [24]. As an example, in social networks, users may be divided into communities based on various factors such as age-group, profession, and racial background. The users' community memberships affects the probability that they are connected with each other. A matching algorithm may potentially use the community membership information to enhance its performance. In order to take the users' community memberships into account, an extension to the CER model is considered which is called the *Stochastic Block Model* (SBM) model. In this model, the edge probabilities depend on their corresponding vertices' community memberships. There have been several works studying the necessary and sufficient conditions for graph matching and the design of practical matching schemes under the SBM model [25]–[28]. However, characterizing tight necessary and sufficient conditions for successful matching and designing polynomial time algorithms which are reliable under these conditions remains an open problem both in the CER and SBM settings.

A further extension of the problem, called '*seeded graph matching*' has also been investigated in the literature [6], [29]–[37]. Seeded graph matching models applications where the matching agent has access to additional side-information in the form of pre-matched *seeds*. A seed vertex is one whose correct label in both graphs is known prior to the start of the matching process. One pertinent application of seeded graph matching is de-anonymization of users over multiple social networks. Many Web users are members of multiple online social networks such as Facebook, Twitter, Google+, LinkedIn, etc. Each online network represents a subset of the users' "real" ego-networks. Graph matching provides algorithms to de-anonymize the users by reconciling these online network graphs, that is, to identify all the accounts belonging to the same individual. In this context, the availability of seeds is justified by the fact that a small fraction of individuals explicitly link their accounts across multiple networks. In this case, these linked accounts can be used as seeds in the matching algorithm. It turns out, that in many cases, these connections may be leveraged to identify a

very large fraction of the users in the network [30]–[34]. In parallel to the study of fundamental limits of graph matching described above, the design of practical low complexity matching algorithms has also been studied in [38]–[40], where reliable matching of real-world networks with up to millions of nodes have been performed.

In this work, we construct an information theoretic framework based on concentration of measure theorems in order to investigate the fundamental limits of graph matching. We propose the '*typicality matching*' (TM) strategy which operates based on the concept of typicality of sequences of random variables [41], and is applicable under a wide range of graph models including CER, SBM and seeded graph matching. The strategy considers the pair of adjacency matrices corresponding to the two graphs. Each $n \times n$ adjacency matrix may be viewed as an n^2 -length sequence of random variables, where n is the number of vertices in the graph. Consequently, one may naturally extend the notion of typicality of sequences of random variables to that of random adjacency matrices. The TM strategy finds a labeling for the vertices in the anonymized graph which results in a pair of jointly typical adjacency matrices for the two graphs, where typicality is defined with respect to the underlying joint edge distribution. The success of the matching algorithm is investigated as the graph size grows asymptotically large. The matching algorithm is said to succeed if the fraction of correctly matched vertices approaches one as the number of vertices goes to infinity. As a result, the TM algorithm is successful as long as any labeling which leads to a pair of jointly typical adjacency matrices assigns an incorrect label to a negligible fraction of size $o(n)$ vertices in the anonymized graph.¹ In order to study the conditions for the success of the TM strategy, we derive several new bounds on the probability of joint typicality of permutations of sequences of random variables. The bounds may be of independent interest in other research areas as well.

The generality of the information theoretic approach allows us to investigate matching under a wide range of statistical models. The contributions of this work can be summarized as follows.

- We build upon the ideas in [7], [25] to develop a general framework based on TM which allows for derivation of necessary and sufficient conditions under which graph matching is possible in a wide range of statistical models. The framework is applicable in matching graphs with weighted edges as well as simultaneous matching of more than two graphs in seeded and seedless matching.
- We apply the TM framework to graph matching under the CER, SBM and seeded graph matching models and to derive theoretical guarantees for successful matching.
- We derive converse results which characterize conditions under which matching is not possible in the CER model as well as simultaneous matching of more than two graphs.
- We investigate the approach proposed in [35], which builds upon the TM framework to propose a polynomial

¹We write $f(x) = o(g(x))$ if $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = 0$.

TABLE I
NOTATION TABLE: RANDOM GRAPHS

g :	unlabeled graph	\mathcal{V} :	vertex set	\mathcal{E} :	edge set
σ :	labeling	\tilde{g} :	labeled graph	G_σ :	adjacency matrix
U_σ :	upper-triangle	ℓ :	# of attributes	\tilde{g}^2 :	relabelled graph

time matching algorithm for the seeded graph matching scenario.

The rest of this article is organized as follows: Section II describes the notation. Section III provides the problem formulation. Section IV develops the necessary tools for analyzing the performance of the TM algorithm. Section V studies matching under the CER model. Section VI considers the SBM model. Section VII investigates matching collections of more than two graphs. In Section VIII, necessary conditions and converse results for matching of pairs of graphs are investigated. Section IX studies seeded graph matching. Section X concludes this article.

II. NOTATION

We represent random variables by capital letters such as X, U and their realizations by small letters such as x, u . Sets are denoted by calligraphic letters such as \mathcal{X}, \mathcal{U} . The set of natural numbers, and the real numbers are represented by \mathbb{N} , and \mathbb{R} respectively. The random variable $\mathbb{1}_{\mathcal{E}}$ is the indicator function of the event \mathcal{E} . The set of numbers $\{n, n+1, \dots, m\}$, $n, m \in \mathbb{N}$ is represented by $[n, m]$. Furthermore, for the interval $[1, m]$, we sometimes use the shorthand notation $[m]$ for brevity. For a given $n \in \mathbb{N}$, the n -length vector (x_1, x_2, \dots, x_n) is written as x^n . We write $a \doteq b \pm \epsilon$ to denote $b - \epsilon \leq a \leq b + \epsilon$, where $a, b, \epsilon \in \mathbb{R}$. We define $|a|^+ \triangleq \max(0, a)$. The notation $\exp_2(\alpha)$ is used to represent 2^α to help readability.

III. PROBLEM FORMULATION

A graph $g = (\mathcal{V}, \mathcal{E})$ is characterized by the vertex set $\mathcal{V} = \{v_1, v_2, \dots, v_n\}$, and the edge set \mathcal{E} . We consider weighted graphs, where each edge is assigned an *attribute* $x \in [0, l-1]$ and $l \geq 2$. Consequently, the edge set \mathcal{E} is a subset of the set $\{(x, v_i, v_j) | i \neq j, x \in [0, l-1]\}$, where for each pair (v_i, v_j) there is a unique attribute x for which $(x, v_i, v_j) \in \mathcal{E}$. For instance, an unlabeled graph with binary valued edges is a graph for which $l = 2$. In this case, if the pair $v_{n,i}$ and $v_{n,j}$ are not connected, we write $(0, v_{n,i}, v_{n,j}) \in \mathcal{E}$, otherwise $(1, v_{n,i}, v_{n,j}) \in \mathcal{E}$. The edge attribute models the nature of the connection between the corresponding vertices. For instance in social network graphs, where vertices represent the members of the network and edges capture their connections, an edge may take different attributes depending on whether the members are family members, close friends, or acquaintances. A labeled graph $\tilde{g} = (g, \sigma)$ is a graph equipped with a bijective *labeling function* $\sigma : \mathcal{V} \rightarrow [n]$. The labeling represents the identity of the members in the social network. For a labeled graph \tilde{g} , the adjacency matrix is defined as $G_\sigma = [g_{\sigma(i,j)}]_{i,j \in [1,n]}$, where $g_{\sigma(i,j)}$ is the unique value such that $(g_{\sigma(i,j)}, v_k, v_l) \in \mathcal{E}$, where $(v_k, v_l) = (\sigma^{-1}(i), \sigma^{-1}(j))$. The

adjacency matrix captures the edge attributes of the graph. The upper triangle (UT) corresponding to \tilde{g} is the structure $U_\sigma = [g_{\sigma(i,j)}]_{i < j}$. The subscript ‘ σ ’ is dropped when there is no ambiguity. The notation is summarized in Table I.

We consider graphs whose edges are generated stochastically. Under the CER and SBM models, we consider special instances of the following random graph model.

Definition 1 (Random Graph): A random graph \tilde{g} generated based on $\prod_{i \in [n], j < i} P_{X_{i,j}}$ is an undirected labeled graph, where the edge between $v_i, i \in [n]$ and $v_j, j < i$ is generated according to $P_{X_{\sigma(i), \sigma(j)}}$ independently of the other edges. Alternatively,

$$P((x, v_i, v_j) \in \mathcal{E}) = P_{X_{\sigma(i), \sigma(j)}}(x), x \in [0, l-1], i, j \in [n].$$

In the graph matching problem, we are given a pair correlated graphs $(\tilde{g}^1, \tilde{g}^2)$, where only the labeling for the vertices of the first graph is available. The objective is to recover the labeling of the vertices in the second graph by leveraging the correlation among their edges. A pair of correlated random graphs is defined below.

Definition 2 (Correlated Random Graph): A pair of correlated random graphs $(\tilde{g}^1, \tilde{g}^2)$ generated based on $\prod_{i \in [n], j < i} P_{X_{i,j}^1, X_{i,j}^2}$ is a pair of undirected labeled graphs. Let v^1, w^1 and v^2, w^2 be two pairs of vertices with the same label in \tilde{g}^1 and \tilde{g}^2 , respectively, i.e., $\sigma^1(v^1) = \sigma^2(v^2) = s_1$ and $\sigma^1(w^1) = \sigma^2(w^2) = s_2$. Then, the pair of edges between (v^1, w^1) and (v^2, w^2) are generated according to $P_{X_{s_1, s_2}^1, X_{s_1, s_2}^2}$. Alternatively,

$$P((x^1, v_i^1, w_j^1) \in \mathcal{E}^1, (x^2, v_i^2, w_j^2) \in \mathcal{E}^2) \\ = P_{X_{s_1, s_2}^1, X_{s_1, s_2}^2}(x^1, x^2), x^1, x^2 \in [0, l-1], i, j \in [n].$$

Remark 1: In Definition 2, the pair $(\tilde{g}^1, \tilde{g}^2)$ are said to be a correlated pair of Erdős-Rényi (CER) graphs if there exists a distribution P_{X^1, X^2} such that $P_{X_{s_1, s_2}^1, X_{s_1, s_2}^2} = P_{X^1, X^2} \forall s_1, s_2 \in [n]$.

A graph matching strategy takes $(\tilde{g}^1, \tilde{g}^2)$ as its input and outputs (\tilde{g}^1, \hat{g}^2) , where \hat{g}^2 is the graph \tilde{g}^2 with its labels σ^2 removed, and \hat{g}^2 is the relabeled graph after matching. The matching strategy is said to succeed if the fraction of correctly matched vertices approaches one as the number of vertices is increased asymptotically. This is formalized below.

Definition 3 (Matching Strategy): For a family of pairs of correlated random graphs $\tilde{g}_n^1 = (g_n^1, \sigma_n^1)$ and $\tilde{g}_n^2 = (g_n^2, \sigma_n^2)$, $n \in \mathbb{N}$, generated based on $\prod_{i \in [n], j < i} P_{X_{i,j}^1, X_{i,j}^2}$, $n \in \mathbb{N}$ where n is the number of vertices. A matching strategy is a sequence of functions $f_n : (\tilde{g}_n^1, \tilde{g}_n^2) \rightarrow (\tilde{g}_n^1, \hat{g}_n^2)$, $n \in \mathbb{N}$, where $\hat{g}_n^2 = (g_n^2, \hat{\sigma}_n^2)$ and $\hat{\sigma}_n^2$ is the reconstruction of σ^2 . Let I_n be

distributed uniformly over $[n]$. The matching strategy is said to succeed if $P(\sigma^2(v_{l_n}^2) = \hat{\sigma}^2(v_{l_n}^2)) \rightarrow 1$ as $n \rightarrow \infty$.

Note that in the above definition, for f_n to succeed, the fraction of vertices whose labels are matched incorrectly must vanish as n approaches infinity. This is a relaxation of the criteria considered in [16]–[21], [35] where all of the vertices are required to be matched correctly simultaneously with vanishing probability of error as $n \rightarrow \infty$. As observed in the next sections, this relaxation leads to a significant simplification in the performance analysis of the proposed matching strategies and allows us to use the concentration of measure theorems and results from information theory to derive theoretical guarantees on the performance of the TM strategy.

Definition 4 (Achievable Region): A family of sets of distributions $\tilde{\mathcal{P}} = (\mathcal{P}_n)_{n \in \mathbb{N}}$ is in the achievable region if for every sequence of distributions $\prod_{s_1 \in [n], s_2 < s_1} P_{X_{s_1, s_2}^1, X_{s_1, s_2}^2}^{(n)} \in \mathcal{P}_n$, there exists a successful matching strategy. The maximal achievable family of sets of distributions is denoted by \mathcal{P}^* .

In social network deanonymization, among other applications, often the correct label of a fraction of the vertices in the anonymized graph are known beforehand. This is due to a fraction of members having used the same user IDs across graphs, or having linked their accounts externally. In these scenarios, the matching strategy may use these pre-matched vertices as ‘seeds’ to recover the labels of the rest of the vertices. Such matching strategies, which are called seeded matching strategies, are defined rigorously and studied in Section IX.

IV. PERMUTATIONS OF TYPICAL SEQUENCES

In the previous section, we described correlated pairs of random graphs, where the graph edges are generated randomly based on an underlying joint distribution. Alternatively, the adjacency matrices of the graphs are generated according to a joint distribution. Furthermore, as explained in Definition 2, we assume that each edge pair connecting two similarly labeled vertices in the two graphs is generated independently of all other edges based on the distribution $P_{X_{i,j}^1, X_{i,j}^2}$, where i, j are the vertex labels. Consequently, it is expected, given large enough graph sizes, that the adjacency matrices of the graphs look ‘typical’ with respect to the joint edge distribution. Roughly speaking, this requires the frequency of joint occurrence of symbols (x^1, x^2) to be close to $\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n P_{X_{i,j}^1, X_{i,j}^2}(x^1, x^2)$, where $x^1, x^2 \in [0, l-1]$. Based on this observation, in the next sections we propose the typicality matching strategy which operates by finding the labeling for the second graph which results in a jointly typical pair of adjacency matrices.

This is analogous to typicality decoding in the channel coding problem in information theory, where the decoder finds the transmitted sequence by searching for a codeword which is jointly typical with the received sequence. In this analogy which is shown in Figure 2, the labeled graph \tilde{g}^2 is passed through a ‘channel’ which outputs the graph g^2 whose labels have undergone a randomly and uniformly chosen permutation, and the matching algorithm acting as a ‘decoder’ wants to recover \tilde{g}^2 using g^2 and the side-information \tilde{g}^1 . Changing the

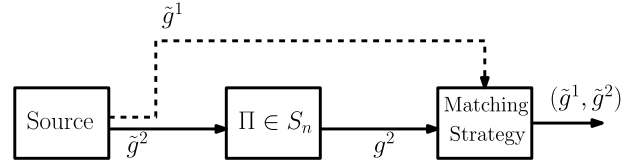


Fig. 2. The pair of correlated graphs $(\tilde{g}^1, \tilde{g}^2)$ are generated as described in Definition 2. The labels in \tilde{g}^2 undergo a random permutation Π chosen uniformly among the set of all possible permutations of n -length sequences S_n . The matching strategy uses \tilde{g}^1 as side information to recover \tilde{g}^2 from g^2 .

labeling of g^2 leads to a permutation of its adjacency matrix. Hence, we need to search over permutations of the adjacency matrix and find the one which leads to a typical pair of adjacency matrices. The error analysis of the TM strategy requires investigating the probability of joint typicality of permutations of pairs of correlated sequences.

In this section, we analyze the joint typicality of permutations of collections of correlated sequences of random variables. While the analysis is used in the subsequent sections to derive the necessary and sufficient conditions for successful matching in various graph matching scenarios, it may also be of independent interest in other research areas as well.

We follow the notation used in [42] in our study of permutation groups summarized below.

Definition 5 (Set Permutation): A permutation on the set of numbers $[1, n]$ is a bijection $\pi : [1, n] \rightarrow [1, n]$. The set of all permutations on the set of numbers $[1, n]$ is denoted by S_n .

Definition 6 (Cycle and Fixed Point): A permutation $\pi \in S_n, n \in \mathbb{N}$ is called a cycle if there exists $k \in [1, n]$ and $\alpha_1, \alpha_2, \dots, \alpha_k \in [1, n]$ such that i) $\pi(\alpha_i) = \alpha_{i+1}, i \in [1, k-1]$, ii) $\pi(\alpha_k) = \alpha_1$, and iii) $\pi(\beta) = \beta$ if $\beta \neq \alpha_i, \forall i \in [1, k]$. The variable k is the length of the cycle. The element β is a fixed point of the permutation if $\pi(\beta) = \beta$. We write $\pi = (\alpha_1, \alpha_2, \dots, \alpha_k)$. The cycle π is non-trivial if $k \geq 2$.

Lemma 1 [42]: Every permutation $\pi \in S_n, n \in \mathbb{N}$ has a unique decomposition into disjoint non-trivial cycles.

Definition 7: For a given $n, m, c \in \mathbb{N}$, and $1 \leq i_1 \leq i_2 \leq \dots \leq i_c \leq n$ such that $n = \sum_{j=1}^c i_j + m$, an $(m, c, i_1, i_2, \dots, i_c)$ -permutation is a permutation in S_n which has m fixed points and c disjoint cycles with lengths i_1, i_2, \dots, i_c , respectively.

Example 1: Consider the permutation which maps the vector $(1, 2, 3, 4, 5)$ to $(5, 1, 4, 3, 2)$. The permutation can be written as a decomposition of disjoint cycles in the following way $\pi = (1, 2, 5)(3, 4)$, where $(1, 2, 5)$ and $(3, 4)$ are cycles with lengths 3 and 2, respectively. The permutation π is a $(0, 2, 2, 3)$ -permutation.

Definition 8 (Sequence Permutation): For a given sequence $y^n \in \mathbb{R}^n$ and permutation $\pi \in S_n$, the sequence $z^n = \pi(y^n)$ is defined as $z^n = (y_{\pi(i)})_{i \in [1, n]}$.²

Definition 9 (Derangement): A permutation on vectors of length n is a derangement if it has no fixed points. The number of distinct derangements of n -length vectors is denoted by $!n$. In our analysis, we make extensive use of the standard permutations defined below.

²Note that in Definitions 5 and 8 we have used π to denote both a scalar function which operates on the set $[1, n]$ as well as a function which operates on the vector space \mathbb{R}^n .

Definition 10 (Standard Permutation): Let $m, c, i_1, i_2, \dots, i_c$ be as in Definition 7. The $(m, c, i_1, i_2, \dots, i_c)$ -standard permutation is defined as the $(m, c, i_1, i_2, \dots, i_c)$ -permutation consisting of the cycles $(\sum_{j=1}^{k-1} i_j + 1, \sum_{j=1}^{k-1} i_j + 2, \dots, \sum_{j=1}^k i_j), k \in [1, c]$. Alternatively, the $(m, c, i_1, i_2, \dots, i_c)$ -standard permutation is defined as:

$$\pi = (1, 2, \dots, i_1)(i_1 + 1, i_1 + 2, \dots, i_1 + i_2) \cdots \left(\sum_{j=1}^{c-1} i_j + 1, \sum_{j=1}^{c-1} i_j + 2, \dots, \sum_{j=1}^c i_j \right) (n - m + 1)(n - m + 2) \cdots (n).$$

Example 2: The $(2, 2, 3, 2)$ -standard permutation is a permutation which has $m = 2$ fixed points and $c = 2$ cycles. The first cycle has length $i_1 = 3$ and the second cycle has length $i_2 = 2$. It is a permutation on sequences of length $n = \sum_{j=1}^c i_j + m = 3 + 2 + 2 = 7$. The permutation is given by $\pi = (123)(45)(6)(7)$. For an arbitrary sequence $\underline{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_7)$, we have:

$$\pi(\underline{\alpha}) = (\alpha_3, \alpha_1, \alpha_2, \alpha_5, \alpha_4, \alpha_6, \alpha_7).$$

A. Typicality of Permutations of Pairs of Correlated Sequences

Definition 11 (Type of Sequences): For a sequence $x^n \in \mathcal{X}^n$, the corresponding type vector $\underline{t} = (t(x))_{x \in \mathcal{X}}$ is defined as $\underline{t}(x) = \frac{\sum_{i=1}^n \mathbb{1}(x_i = x)}{n}, x \in \mathcal{X}$. For the pair $(x^n, y^n) \in \mathcal{X}^n \times \mathcal{Y}^n$, the corresponding joint type $\underline{s} = (s(x, y))_{x, y \in \mathcal{X} \times \mathcal{Y}}$ is defined as $\underline{s}(x, y) = \frac{\sum_{i=1}^n \mathbb{1}(x_i = x, y_i = y)}{n}, x, y \in \mathcal{X} \times \mathcal{Y}$.

Definition 12 (Strong Typicality) [41]: Let the pair of random variables (X, Y) be defined on the probability space $(\mathcal{X} \times \mathcal{Y}, P_{X,Y})$, where \mathcal{X} and \mathcal{Y} are finite alphabets. The ϵ -typical set of sequences of length n with respect to $P_{X,Y}$ is defined as:

$$\mathcal{A}_\epsilon^n(X, Y) = \left\{ (x^n, y^n) : \underline{t}(x, y) \doteq P_{X,Y}(x, y) \pm \epsilon, \quad \forall (x, y) \in \mathcal{X} \times \mathcal{Y} \text{ \& } \underline{t}(x, y) = 0 \text{ if } P_{X,Y}(x, y) = 0 \right\},$$

where \underline{t} is the joint type of (x^n, y^n) , $\epsilon > 0$, and $n \in \mathbb{N}$.

For a correlated pair of independent and identically distributed (i.i.d) sequences (X^n, Y^n) and an arbitrary permutation $\pi \in \mathcal{S}_n$, we are interested in bounding the probability $P((X^n, \pi(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y))$. The following proposition shows that in order to find bounds on the probability of joint typicality of permutations of correlated sequences, it suffices to study standard permutations.

Proposition 1: Let (X^n, Y^n) be a pair of i.i.d sequences defined on finite alphabets. We have:

i) For an arbitrary permutation $\pi \in \mathcal{S}_n$,

$$P((\pi(X^n), \pi(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) = P((X^n, Y^n) \in \mathcal{A}_\epsilon^n(X, Y)).$$

ii) Following the notation in Definition 10, let π_1 be an arbitrary $(m, c, i_1, i_2, \dots, i_c)$ -permutation and let π_2 be the $(m, c, i_1, i_2, \dots, i_c)$ -standard permutation. Then,

$$P((X^n, \pi_1(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) = P((X^n, \pi_2(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)).$$

iii) For arbitrary permutations $\pi_x, \pi_y \in \mathcal{S}_n$, let π be the standard permutation having the same number of cycles and cycle lengths as that of $\pi_x^{-1}(\pi_y)$. Then,

$$P((\pi_x(X^n), \pi_y(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) = P((X^n, \pi(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)).$$

iv) For an arbitrary permutation $\pi \in \mathcal{S}_n$,

$$P((X^n, \pi(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) = P((X^n, \pi^{-1}(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)).$$

Proof: Part i) follows from the fact that permuting both X^n and Y^n by the same permutation does not change their joint type. For part ii), it is known that there exists a permutation π such that $\pi(\pi_1) = \pi_2(\pi)$ [42]. Then the statement is proved using part i) as follows:

$$\begin{aligned} P((X^n, \pi_1(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) &= P((\pi(X^n), \pi(\pi_1(Y^n))) \in \mathcal{A}_\epsilon^n(X, Y)) \\ &= P((\pi(X^n), \pi_2(\pi(Y^n))) \in \mathcal{A}_\epsilon^n(X, Y)) \\ &= P((\tilde{X}^n, \pi_2(\tilde{Y}^n)) \in \mathcal{A}_\epsilon^n(X, Y)) \\ &= P((X^n, \pi_2(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)), \end{aligned} \quad (1)$$

where in (1) we have defined $(\tilde{X}^n, \tilde{Y}^n) = (\pi(X^n), \pi(Y^n))$, and (2) holds since $(\tilde{X}^n, \tilde{Y}^n)$ has the same distribution as (X^n, Y^n) . Part iii) follows directly from Parts i) and ii). Part iv) follows from Part ii) by noting that the number and lengths of cycles in π^{-1} is the same as that of π . ■

The following theorem provides an upper-bound on the probability of joint typicality of permutations of correlated sequences for an arbitrary permutation with $m \in [n]$ fixed points.

Theorem 1: Let $\epsilon \in [0, \frac{1}{2} \min_{x,y \in \mathcal{X} \times \mathcal{Y}} P_{X,Y}(x, y)]$, and consider (X^n, Y^n) a pair of i.i.d sequences defined on finite alphabets \mathcal{X} and \mathcal{Y} , respectively. For any permutation π with $m \in [n]$ fixed points, the following holds:

$$\begin{aligned} P((X^n, \pi(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) &\leq 2^{-n(E_\alpha - \zeta_n - \delta_\epsilon)}, \\ E_\alpha &= \min_{t'_X \in \mathcal{P}} \frac{1}{2} \left((1 - \alpha) D(t'_X || P_X) + \alpha D(t'_X || P_X) \right. \\ &\quad \left. + D(P_{X,Y} || (1 - \alpha) P_X P_{Y''} + \alpha P_{X,Y}) \right), \end{aligned} \quad (3)$$

where $\alpha \triangleq \frac{m}{n}$, $\mathcal{P} \triangleq \{t_X \in \mathcal{P}_X | \forall x \in \mathcal{X} : t_X(x) \in \frac{1}{1-\alpha} [P_X(x) - \alpha, P_X(x)]\}$, \mathcal{P}_X is the probability simplex on the alphabet \mathcal{X} , $D(\cdot || \cdot)$ is the Kullback-Leibler divergence, $t'_X \triangleq \frac{1}{\alpha} (P_X - (1 - \alpha) t'_X)$, $P_{Y''}(\cdot) \triangleq \sum_{x \in \mathcal{X}} t'_X(x) P_{Y|X}(\cdot | x)$, $\zeta_n \triangleq \frac{3}{2} |\mathcal{X}|^2 |\mathcal{Y}| \frac{\log(n+1)}{n} + 6 |\mathcal{X}| |\mathcal{Y}| \frac{\log(n+1)}{n}$, and

$$\begin{aligned} \delta_\epsilon &\triangleq |\mathcal{X}| |\mathcal{Y}| \\ &\times \left| \max_{\substack{x,y \in \mathcal{X} \times \mathcal{Y} \\ P_{X,Y}(x,y) \neq 0}} \log \frac{P_{X,Y}(x,y)}{\alpha P_{X,Y}(x,y) + (1-\alpha) P_X(x) P_Y(y)} \right| + O(\epsilon). \end{aligned}$$

The proof is provided in the supplementary material. In the following, we describe an outline of the proof. Let us define \mathcal{A} as the set of fixed points of the permutation π . From the theorem statement, the set \mathcal{A} includes $\alpha = \frac{m}{n}$ fraction of the indices $[n]$. Let \underline{T}_X be the type of the vector of $X_i, i \in \mathcal{A}$, and let \underline{T}_X be the type of the vector $X_i, i \notin \mathcal{A}$. A necessary condition for $(X^n, \pi(Y^n))$ to be jointly ϵ -typical with respect to $P_{X,Y}$

is that $\underline{T}'_X = \underline{t}'_X$ and $\underline{T}''_X = \underline{t}''_X$ such that $(1-\alpha)\underline{t}'_X + \alpha\underline{t}''_X \doteq P_X \pm \epsilon$. Since X^n is an i.i.d sequence of variables, from standard information theoretic arguments, the probability of the event that $\underline{T}'_X = \underline{t}'_X$ decays exponentially in n with exponent $\alpha D(\underline{t}'_X || P_X)$. Similarly, the probability that $\underline{T}''_X = \underline{t}''_X$ decays exponentially in n with exponent $(1-\alpha)D(\underline{t}''_X || P_X)$. This justifies the term $(1-\alpha)D(\underline{t}'_X || P_X) + \alpha D(\underline{t}''_X || P_X)$ in the exponent E_α in Equation (4). Next, note that for $i \in \mathcal{A}$, we have $\pi(Y_i) = Y_i$. As a result, the joint distribution of each of the pairs $(X_i, \pi(Y_i))$, $i \in \mathcal{A}$ is $P_{X,Y}$. On the other hand, for indices $i \notin \mathcal{A}$, we have $\pi(Y_i) = Y_{\pi(i)}$, where $\pi(i) \neq i$. So, the pair $(X_i, \pi(Y_i))$, $i \notin \mathcal{A}$ is an independent pair of variables, where X_i is generated based on P_X , and $\pi(Y_i)$ is generated based on $P_{Y|X}(\cdot | X_{\pi(i)})$. Note that given that $\underline{T}'_X = \underline{t}'_X$, the average distribution of $\pi(Y_i)$, $i \notin \mathcal{A}$ is $\frac{1}{n-|\mathcal{A}|} \sum_{i \notin \mathcal{A}} P_{Y_{\pi(i)}} = \sum_{x \in \mathcal{X}} \underline{t}'_X(x) P_{Y|X}(\cdot | x) = P_{Y''}$. Consequently, the average distribution of $(X_i, \pi(Y_i))$, $i \notin \mathcal{A}$ is $P_X P_{Y''}$. As a result, the average distribution of $(X^n, \pi(Y^n))$ is $(1-\alpha)P_X P_{Y''} + \alpha P_{X,Y}$. Hence, using standard information theoretic arguments, the probability that the pair $(X^n, \pi(Y^n))$ is jointly ϵ -typical with respect to $P_{X,Y}$ decays exponentially with exponent $D(P_{X,Y} || (1-\alpha)P_X P_{Y''} + \alpha P_{X,Y})$, which appears as the third term in the exponent E_α in Equation (4). The exponent E_α can be further simplified for special classes of permutations. For instance, if π does not have any fixed points, the following corollary to Theorem 1 holds.

Corollary 1: If the permutation π in Theorem 1 has no fixed points (i.e., $\alpha = 0$), then:

$$P((X^n, \pi(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) \leq 2^{-n(E_0 - \zeta_n - \delta_\epsilon)}, \quad (5)$$

where, $E_0 = \frac{1}{2}I(X; Y)$, and ζ_n and δ_ϵ are defined in Theorem 1.

Proof: The proof follows from Theorem 1. Note that when $\alpha = 0$, the set \mathcal{P} in the theorem statement has a single element P_X . So, we have $\underline{t}'_X = P_X$, $P_{Y''} = P_Y$, $E_0 = \frac{1}{2}I(X; Y)$, and $\delta_\epsilon = \epsilon |\mathcal{X}||\mathcal{Y}| \left| \max_{x,y \in \mathcal{X} \times \mathcal{Y}: P_{X,Y}(x,y) \neq 0} \log \frac{P_{X,Y}(x,y)}{P_X(x)P_Y(y)} \right|$.

Corollary 1 shows that for two sequences (X^n, Y^n) generated jointly according to $P_{X,Y}$ and a permutation π without any fixed points, the probability that the pair $(X^n, \pi(Y^n))$ are jointly typical with respect to $P_{X,Y}$ decays exponentially in n with exponent $\frac{1}{2}I(X; Y)$. Note that since there are no fixed points in the permutation, each pair $(X_i, \pi(Y_i))$, $i \in [n]$ has joint distribution $P_X P_Y$. So, there is no 'single-letter' correlation among the elements of $(X^n, \pi(Y^n))$. It is well-known that if the sequences X^n and Y^n are generated independently of each other according to marginal distributions P_X and P_Y , respectively, then the probability that they are jointly typical with respect to the joint distribution $P_{X,Y}$ decays exponentially in n with exponent $I(X; Y)$ (e.g., [41]). The coefficient $\frac{1}{2}$ in exponent in Corollary 1 is an artifact of the n -letter correlation between $(X^n, \pi(Y^n))$ due to the fact that the original pair of sequences (X^n, Y^n) are correlated.

Theorem 1 is used in the next sections to derive sufficient conditions under which pairs of correlated graphs can be matched successfully. However, the arguments in the proof of the theorem do not extend naturally to typicality of collections

of more than two permuted sequences. Bounds on the probability of joint typicality of such collections are necessary for evaluating graph matching for collections of graphs studied in Section VII. To this end, the following theorem and the ensuing corollary provide an alternative bound on the probability of joint typicality of $(X^n, \pi(Y^n))$ for an arbitrary permutation π which is then extended to evaluate the typicality of collections of more than two sequences in Theorem 3. This is used in Section VII to evaluate matching of more than two graphs. Additionally, we will observe in the proof of Theorem 5 in Section V that E'_α derived below yields tighter bounds on the probability of joint typicality of $(X^n, \pi(Y^n))$ for large α compared to E_α derived in Theorem 1.

Theorem 2: Let (X^n, Y^n) be a pair of i.i.d sequences defined on finite alphabets \mathcal{X} and \mathcal{Y} , respectively. For any permutation π with $m \in [n]$ fixed points, the following holds:

$$P((X^n, \pi(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) \leq 2^{-n(E'_\alpha - \zeta'_n - \delta_\epsilon)}, \quad (6)$$

$$E'_\alpha = \min_{\underline{t}'_{X,Y} \in \mathcal{P}'} \left(\frac{1-\alpha}{3} \right) D(\underline{t}'_{X,Y} || P_X P_Y) + \alpha D(\underline{t}'_{X,Y} || P_{X,Y}), \quad (7)$$

where $\alpha \triangleq \frac{m}{n}$, $\mathcal{P}' \triangleq \{\underline{t}_{X,Y} \in \mathcal{P}_{X,Y} | \forall (x,y) \in \mathcal{X} \times \mathcal{Y} : \underline{t}_{X,Y}(x,y) \in \frac{1}{1-\alpha}[P_{X,Y}(x,y) - \alpha, P_{X,Y}(x,y)]\}$, $\mathcal{P}_{X,Y}$ is the probability simplex on the alphabet $\mathcal{X} \times \mathcal{Y}$, $\underline{t}'_{X,Y} \triangleq \frac{1}{\alpha}(P_{X,Y} - (1-\alpha)\underline{t}_{X,Y})$, $\zeta'_n = 4|\mathcal{X}||\mathcal{Y}| \log \frac{n+1}{n}$, and δ_ϵ is defined as in Theorem 1.

Proof: In the supplementary material. ■

Computing E'_α requires optimizing over $\underline{t}'_{X,Y}$, which is computationally challenging for large alphabets. The following removes the optimization and provides a lower bound for E'_α .

Corollary 2: Let (X^n, Y^n) be a pair of i.i.d sequences defined on finite alphabets \mathcal{X} and \mathcal{Y} , respectively. For any permutation π with $m \in [n]$ fixed points, the following holds:

$$P((X^n, \pi(Y^n)) \in \mathcal{A}_\epsilon^n(X, Y)) \leq 2^{-n(\widehat{E}'_\alpha - \zeta'_n - \frac{\delta_\epsilon}{3})}, \quad (8)$$

$$\widehat{E}'_\alpha = \frac{1}{3}D(P_{X,Y} || (1-\alpha)P_X P_Y + \alpha P_{X,Y}), \quad (9)$$

where $\alpha \triangleq \frac{m}{n}$, $\zeta'_n \triangleq 4|\mathcal{X}||\mathcal{Y}| \log \frac{n+1}{n}$ and δ_ϵ is defined as in Theorem 1.

Proof: In the supplementary material. ■

It is desirable to find the largest exponent which can be used to bound the exponential decay in the probability of joint typicality of $(X^n, \pi(Y^n))$. Consequently, a question of interest is whether one of the two exponents E_α and E'_α is strictly larger than the other. Towards such a comparison, the next lemma shows that the relation $\frac{2}{3}E_\alpha \leq \widehat{E}'_\alpha \leq E'_\alpha$ holds. On the other hand, it can be shown through analytical evaluations of the bounds under specific distributions $P_{X,Y}$ that in some instances, the relation $E'_\alpha < E_\alpha$ holds. We will observe in the proof of Theorem 5 in Section V that E_α in Theorem 1 yields tighter bounds on the probability of joint typicality when α is small, whereas E'_α in Theorem 2 is useful when evaluating permutations with large α .

Lemma 2: For the exponents E_α , E'_α , and \widehat{E}'_α in Theorems 1 and 2 and Corollary 2, we have:

$$\frac{2}{3}E_\alpha \leq \widehat{E}'_\alpha \leq E'_\alpha.$$

Proof: The relation $\widehat{E}'_\alpha \leq E'_\alpha$ follows by convexity of KL divergence. Also, note that

$$\begin{aligned} \frac{2}{3}E_\alpha &= \frac{2}{3} \min_{t'_X \in \mathcal{P}} \frac{1}{2} \left((1-\alpha)D(t'_X||P_X) + \alpha D(t''_X||P_X) \right. \\ &\quad \left. + D(P_{X,Y}|| (1-\alpha)P_X P_{Y''} + \alpha P_{X,Y}) \right) \\ &\leq \frac{1}{3}D(P_{X,Y}|| (1-\alpha)P_X P_Y + \alpha P_{X,Y}) = \widehat{E}'_\alpha, \end{aligned}$$

where the inequality follows by taking $t'_X = P_X$. Note that this leads to $t''_X = P_X$, so that $D(t'_X||P_X) = D(t''_X||P_X) = 0$. ■

We have provided bounds on the probability of joint typicality of X^n and $\pi(Y^n)$ as a function of the number of fixed points m of the permutation $\pi(\cdot)$. Such bounds are often used in error analysis and derivation of error bounds in various applications [7], [43], [44], and we will use them in the following sections to evaluate the probability of error in graph matching scenarios. In order to evaluate the error exponents, the following results on the limiting behavior of the number of distinct permutations with a given number of fixed points are needed.

Lemma 3: Let $n \in \mathbb{N}$. Let N_m be the number of distinct permutations with exactly $m \in [0, n]$ fixed points. Then,

$$\frac{n!}{m!(n-m)} \leq N_m = \binom{n}{m}!(n-m) \leq n^{n-m}. \quad (10)$$

Particularly, let $m = \alpha n, 0 < \alpha < 1$. Then, the following holds:

$$\lim_{n \rightarrow \infty} \frac{\log N_m}{n \log n} = 1 - \alpha. \quad (11)$$

Proof: In the supplementary material. ■

B. Typicality of Permutations of Collections of Correlated Sequences

We consider joint typicality of permutations of more than two correlated sequences $(X^n_{(1)}, X^n_{(2)}, \dots, X^n_{(k)}), n \in \mathbb{N}, k > 2$. The derivations in this section are used in Section VII to extend the analysis of the TM strategy to simultaneous matching of collections of more than two graphs.

Definition 13 (Strong Typicality of Collections of Sequences [41]): Let the random vector X^k be defined on the probability space $(\prod_{j \in [k]} \mathcal{X}_j, P_{X^k})$, where $\mathcal{X}_j, j \in [k]$ are finite alphabets, and $k > 2$. The ϵ -typical set of sequences of length n with respect to P_{X^k} is defined as:

$$\mathcal{A}_\epsilon^n(X^k) = \left\{ \left(x^n_{(j)} \right)_{j \in [k]} : \underline{t}(\alpha^k) \doteq P_{X^k}(\alpha^k) \pm \epsilon \ \forall \alpha^k \in \prod_{j \in [k]} \mathcal{X}_j \right\},$$

where $\epsilon > 0$, and $\underline{t}(\alpha^k) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}((x_{(j),i})_{j \in [k]} = \alpha^k)$ is the type of $(x^n_{(j)})_{j \in [k]}$.

In the previous section, in order to investigate the typicality of permutations of pairs of correlated sequences, we introduced standard permutations which are completely characterized by the number of fixed points, number of cycles, and cycle lengths of the permutation. The concept of standard permutations does not extend naturally when there are more than two sequences (i.e., more than one non-trivial permutation).

Consequently, investigating typicality of permutations of collections of sequences requires developing additional analytical tools described next.

Definition 14 (Bell Number [45]): Let $\mathbf{P} = \{\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_{b_k}\}$ be the set of all partitions of $[1, k]$. The natural number b_k is the k 'th Bell number. We make the convention that $\mathcal{P}_{b_k} = \{[n]\}$.

In the following, we define Bell permutation vectors which are analogous to standard permutations for the case when the problem involves more than one non-trivial permutation.

Definition 15 (Partition Correspondence): Let $k, n \in \mathbb{N}$ and $(\pi_1, \pi_2, \dots, \pi_k)$ be arbitrary permutations operating on n -length vectors. The index $i \in [1, n]$ is said to correspond to the partition $\mathcal{P}_j \in \mathbf{P}$ of the set $[1, k]$ if the following holds:

$$\forall l, l' \in [1, k] : \pi_l^{-1}(i) = \pi_{l'}^{-1}(i) \iff \exists r : l, l' \in \mathcal{D}_{j,r},$$

where $\mathcal{P}_j = \{\mathcal{D}_{j,1}, \mathcal{D}_{j,2}, \dots, \mathcal{D}_{j,|\mathcal{P}_j|}\}$.

Example 3: Let us consider a triple of permutations of n -length sequences, i.e., $k = 3$, and the partition $\mathcal{P} = \{\{1, 2\}, \{3\}\}$. An index $i \in [n]$ corresponds to \mathcal{P} if the first two permutations map the index to the same integer and the third permutation maps the index to a different integer.

Definition 16 (Bell Permutation Vector): Let $(i_1, i_2, \dots, i_{b_k})$ be an arbitrary sequence, where $\sum_{k \in [b_k]} i_k = n, i_k \in [0, n]$, b_k is the k th Bell number, and $n, k \in \mathbb{N}$. The vector of permutations $(\pi_1, \pi_2, \dots, \pi_k)$ is called an $(i_1, i_2, \dots, i_{b_k})$ -Bell permutation vector if for every partition \mathcal{P}_k exactly i_k indices correspond to that partition. Equivalently:

$$\begin{aligned} \forall j \in [b_k] : i_k = \left| \left\{ i \in [n] : \forall l, l' \in [k] : \pi_l^{-1}(i) = \pi_{l'}^{-1}(i) \right. \right. \\ \left. \left. \iff \exists r \in [|\mathcal{P}_j|] : l, l' \in \mathcal{D}_{j,r} \right\} \right|, \end{aligned}$$

where $\mathcal{P}_j = \{\mathcal{D}_{j,1}, \mathcal{D}_{j,2}, \dots, \mathcal{D}_{j,|\mathcal{P}_j|}\}$.

The definition of Bell permutation vectors is further clarified through the following example.

Example 4: Consider three permutations (π_1, π_2, π_3) of vectors with length seven, i.e., $k = 3$ and $n = 7$. Then, $b_k = 5$ and we have:

$$\begin{aligned} \mathcal{P}_1 &= \{\{1\}, \{2\}, \{3\}\}, \quad \mathcal{P}_2 = \{\{1, 2\}, \{3\}\}, \quad \mathcal{P}_3 = \{\{1, 3\}, \{2\}\}, \\ \mathcal{P}_4 &= \{\{1\}, \{2, 3\}\}, \quad \mathcal{P}_5 = \{\{1, 2, 3\}\}. \end{aligned}$$

Let π_1 be the identity permutation, $\pi_2 = (135)(24)$, and $\pi_3 = (15)(24)(37)$. Then:

$$\begin{aligned} \pi_1((1, 2, \dots, 7)) &= (1, 2, 3, 4, 5, 6, 7), \\ \pi_2((1, 2, \dots, 7)) &= (5, 4, 1, 2, 3, 6, 7), \\ \pi_3((1, 2, \dots, 7)) &= (5, 4, 7, 2, 1, 6, 3), \end{aligned}$$

The vector (π_1, π_2, π_3) is a $(2, 1, 0, 3, 1)$ -Bell permutation vector, where the indices $(3, 5)$ correspond to the \mathcal{P}_1 partition (each of the three permutations map indices $(3, 5)$ to a different integer), index 7 corresponds to the \mathcal{P}_2 partition (the first two permutations map the index 7 to the same integer which is different from the one for the third permutation), indices $(1, 2, 4)$ correspond to the \mathcal{P}_4 permutation (the second and third permutations map the indices $(1, 2, 4)$ to the same integer which is different from the output of the first permutation), and index 6 corresponds to \mathcal{P}_5 (all permutations map

the index 6 to the same integer). None of the indices corresponds to \mathcal{P}_3 since there is no index which is mapped to the same integer by the first and third permutations and a different integer by the second permutation.

Remark 2: Bell permutation vectors are not unique. There can be several distinct $(i_1, i_2, \dots, i_{b_k})$ -Bell permutation vectors for given $n, k, i_1, i_2, \dots, i_{b_k}$. This is in contrast with standard permutations defined in Definition 10, which are unique given the parameters $n, k, c, i_1, i_2, \dots, i_c$.

The following theorem provides bounds on the probability of joint typicality of permutations of collections of correlated sequences.

Theorem 3: Let $(X_{(j)}^n)_{j \in [k]}$ be a collection of correlated sequences of i.i.d random variables defined on finite alphabets $\mathcal{X}_j, j \in [k]$. For any $(i_1, i_2, \dots, i_{b_k})$ -Bell permutation vector $(\pi_1, \pi_2, \dots, \pi_k)$, the following holds:

$$P\left(\left(\pi_i(X_{(j)}^n)_{j \in [k]} \in \mathcal{A}_\epsilon^n(X^k)\right)\right) \leq 2^{-n(E_{i_1, i_2, \dots, i_{b_k}} + O(\epsilon) + O(\frac{\log n}{m}))}, \quad (12)$$

$$E_{i_1, i_2, \dots, i_{b_k}} = -\frac{1}{(k(k-1)+1)(b_k-1)} D\left(P_{X^k} \parallel \sum_{j \in [b_k]} \frac{i_j}{n} P_{X_{\mathcal{P}_j}}\right) \quad (13)$$

where $P_{X_{\mathcal{P}_j}} = \prod_{r \in [1, |\mathcal{P}_j|]} P_{X_{i_1, i_2, \dots, i_{|\mathcal{D}_{j,r}|}}}$, $\mathcal{D}_{j,r} = \{i_1, i_2, \dots, i_{|\mathcal{D}_{j,r}|}\}, j \in [b_k], r \in [1, |\mathcal{P}_j|]$.

Proof: In the supplementary material. ■

Note that for the special case of permutations of pairs of sequences of random variables, $k = 2$, the second Bell number is $b_k = 2$. In this case $(k(k-1)+1)(b_k-1) = 3$, and the bound on the probability of joint typicality given in Theorem 3 is the same as the one in Corollary 2.

In the following, we generalize Lemma 3 to the case where a collection of more than two permuted sequence is considered, and provide upper and lower bounds on the number of distinct Bell permutation vectors for a given vector $(i_1, i_2, \dots, i_{b_k})$.

Definition 17 (k-fold Derangement): A vector $(\pi_1(\cdot), \pi_2(\cdot), \dots, \pi_k(\cdot))$ of permutations of n -length sequences is called an k -fold derangement if $\pi_l(\cdot)$ is the identity permutation, and $\pi_l(i) \neq \pi_{l'}(i), l, l' \in [k], l \neq l', i \in [n]$. The number of distinct k -fold derangements of $[n]$ is denoted by $d_k(n)$. Particularly $d_2(n) = !n$ is the number of derangements of $[n]$.

Lemma 4: Let $n \in \mathbb{N}$ and $k \in [n]$. Then,

$$((n-k+1)!)^{k-1} \leq d_k(n) \leq (!n)^{k-1}.$$

Proof: In the supplementary material. ■

Lemma 5: Let $(i_1, i_2, \dots, i_{b_k})$ be a vector of non-negative integers such that $\sum_{j \in [b_k]} i_j = n$. Define $N_{i_1, i_2, \dots, i_{b_k}}$ as the number of distinct $(i_1, i_2, \dots, i_{b_k})$ -Bell permutation vectors. Then,

$$\begin{aligned} \binom{n}{i_1, i_2, \dots, i_{b_k}} \prod_{j \in [b_k]} d_{|\mathcal{P}_j|}(i_j) &\leq N_{i_1, i_2, \dots, i_{b_k}} \\ &\leq \binom{n}{i_1, i_2, \dots, i_{b_k}} n^{\sum_{j \in [b_k]} |\mathcal{P}_j| i_j - n}. \end{aligned} \quad (14)$$

Particularly, let $i_k = \alpha_k \cdot n, n \in \mathbb{N}$. The following holds:

$$\lim_{n \rightarrow \infty} \frac{\log N_{i_1, i_2, \dots, i_{b_k}}}{n \log n} = \sum_{j \in [b_k]} |\mathcal{P}_j| \alpha_j - 1. \quad (15)$$

Proof: In the supplementary material. ■

V. MATCHING ERDŐS-RÉNYI GRAPHS

In this section, we consider matching of CPER graphs with weighted edges. In Section III, we described correlated random graphs. A CPER is a special instance of the correlated random graphs defined in Definition 2. We propose the typicality matching strategy and provide sufficient conditions on the joint edge statistics under which the strategy succeeds.

A. The Typicality Matching Strategy for CPERs

Given a correlated pair of graphs (\tilde{g}^1, g^2) , where only the labeling for \tilde{g}^1 is given, the TM strategy operates as follows. The scheme finds a labeling $\hat{\sigma}^2$, for which the pair of UT's $U_{\sigma^1}^1$ and $U_{\hat{\sigma}^2}^2$ are jointly typical with respect to $P_{X_1, X_2}^{(n)}$ when viewed as vectors of length $\frac{n(n-1)}{2}$. The strategy fails if no such labeling exists. Alternatively, it finds an element $\hat{\sigma}^2$ in the set:

$$\hat{\Sigma} = \left\{ \hat{\sigma}^2 \mid (U_{\sigma^1}^1, U_{\hat{\sigma}^2}^2) \in \mathcal{A}_\epsilon^{\frac{n(n-1)}{2}}(X_1, X_2) \right\}, \quad (16)$$

where $\epsilon = \omega(\frac{1}{n})$. The algorithm declares $\hat{\sigma}^2$ as the correct labeling. Note that the set $\hat{\Sigma}$ may have more than one element. In that case, the strategy chooses one of these elements as the output randomly. We will show that under certain conditions on the joint graph statistics, all of the elements of $\hat{\Sigma}$ satisfy the criteria for successful matching given in Definition 3. In other words, for all of the elements of $\hat{\Sigma}$ the probability of incorrect labeling for any given vertex is arbitrarily small for large n . Formally, The TM strategy is a sequence of functions $f_n : (\tilde{g}_n^1, g_n^2) \rightarrow (\tilde{g}_n^1, \hat{g}_n^2), n \in \mathbb{N}$, where for any given $n \in \mathbb{N}$, the labeling $\hat{\sigma}_n^2$ of g_n^2 is chosen randomly and uniformly from the set $\hat{\Sigma}$ defined in Equation (16).

Theorem 4: For the TM strategy described above, a given family of sets of distributions $\tilde{\mathcal{P}} = (\mathcal{P}_n)_{n \in \mathbb{N}}$ is achievable, if for every sequence of distributions $P_{X_1, X_2}^{(n)} \in \mathcal{P}_n, n \in \mathbb{N}$,

$$2(1-\alpha) \frac{\log n}{n-1} \leq \max(E_{\alpha^2}, E'_{\alpha^2}), 0 \leq \alpha \leq \alpha_n, \quad (17)$$

and $\max_{(x_1, x_2) : P_{X_1, X_2}^{(n)}(x_1, x_2) \neq 0} \left| \log \frac{P_{X_1}^{(n)}(x_1) P_{X_2}^{(n)}(x_2)}{P_{X_1, X_2}^{(n)}(x_1, x_2)} \right|^+ = o(\log n)$, where $\alpha_n \rightarrow 1$ as $n \rightarrow \infty$, and E_{α^2} and E'_{α^2} are defined in Theorems 1 and 2, respectively.

Proof: In the supplementary material. ■

Theorem 4 provides sufficient conditions on the edge statistics of the CPER graphs such that the TM strategy correctly matches ‘almost’ all vertices of the two graphs. That is, the theorem provides sufficient conditions under which $P(\sigma^2(v_n^2) = \hat{\sigma}^2(v_n^2)) \rightarrow 1$ as $n \rightarrow \infty$, where I_n is chosen uniformly among all indices $[n]$, as defined in Definition 3. A question of significant interest is how this sufficient condition changes if the success criterion is relaxed

so that the strategy is only required to correctly match a fraction $\beta \in [0, 1]$ of the vertices. More precisely, we want to know the conditions on $P_{X_1, X_2}^{(n)}, n \in \mathbb{N}$ such that $P\left(\frac{1}{n} \left| \left\{ i : \sigma^2(v_i^2) = \hat{\sigma}^2(v_i^2) \right\} \right| \geq \beta \right) \rightarrow 1$ as $n \rightarrow \infty$. This is of particular interest in social network deanonymization [5]–[7], where even a small fraction of matched vertices is a violation of those users' privacy. The following corollary to Theorem 4 provides sufficient conditions under which such partial matching is possible using the TM strategy.

Corollary 3 (Partial Matching of CPERs): Let $\beta > 0$. Given a sequence of distributions $P_{X_1, X_2}^{(n)}$, the TM strategy correctly matches at least β fraction of the vertices for asymptotically large n (i.e., $P\left(\frac{1}{n} \left| \left\{ i : \sigma^2(v_i^2) = \hat{\sigma}^2(v_i^2) \right\} \right| \geq \beta \right) \rightarrow 1$ as $n \rightarrow \infty$), given that the following holds:

$$2(1 - \alpha) \frac{\log n}{n - 1} \leq \max(E_{\alpha^2}, E'_{\alpha^2}), 0 \leq \alpha \leq \beta, \quad (18)$$

and $\max_{(x_1, x_2): P_{X_1, X_2}^{(n)}(x_1, x_2) \neq 0} \left| \log \frac{P_{X_1}(x_1)P_{X_2}(x_2)}{P_{X_1, X_2}^{(n)}(x_1, x_2)} \right|^+ = o(\log n)$.

The proof follows from the proof of Theorem 4 in the supplementary material by replacing α_n with β .

Remark 3: We have restricted our analysis to matching undirected graphs where $(x, v_i, v_j) \in \mathcal{E}$ if and only if $(x, v_j, v_i) \in \mathcal{E}$. The results can be extended to directed graphs by evaluating the joint typicality of the complete adjacency matrices of the graphs rather than the upper-triangles of the adjacency matrices.

B. Matching Under the Erasure Model

In the following, we consider matching pairs of CERs under the special case of the erasure model, where the following distribution on the graph edges is considered:

$$\begin{aligned} P_{X_1, X_2}^{(n)}(0, 0) &= 1 - p_n, & P_{X_1, X_2}^{(n)}(0, 1) &= 0, \\ P_{X_1, X_2}^{(n)}(1, 0) &= p_n(1 - s), & P_{X_1, X_2}^{(n)}(1, 1) &= p_ns, \\ P_X^{(n)}(0) &= 1 - p_n, & P_X^{(n)}(1) &= p_n, \\ P_Y^{(n)}(0) &= 1 - p_ns, & P_Y^{(n)}(1) &= p_ns \end{aligned}$$

where $s \in [0, 1]$ is fixed in n , and $p_n \rightarrow 0$ as $n \rightarrow \infty$. The model has been studied extensively in the literature (e.g., [17], [29], [30]). Under the erasure model, there is an edge between each two vertices in \tilde{g}_n^1 with probability p_n . The edges in \tilde{g}_n^2 are sampled from the edges in \tilde{g}_n^1 such that each edge in \tilde{g}_n^1 is erased in \tilde{g}_n^2 with probability $(1 - s)$ and it is kept with probability s . We are interested in finding the fastest rate at which $p_n \rightarrow 0$ such that successful matching is possible. In [30], it is shown that a sufficient condition for successful matching under the erasure model is that $\frac{s}{2}p_n \geq \frac{\ln n}{n}$ as $n \rightarrow \infty$, where \ln is the natural logarithm. Alternatively, a sufficient condition for successful matching is $\lim_{n \rightarrow \infty} \frac{\frac{\ln n}{n}}{p_n} \leq \frac{s}{2}$. The following theorem shows that the TM strategy improves this sufficient condition for successful matching.

Theorem 5: Let $\frac{1}{4} < s < \frac{1}{2}$. There exists a sequence p_n approaching 0 as $n \rightarrow \infty$ for which i) the TM strategy leads to successful matching, and ii) $\lim_{n \rightarrow \infty} \frac{\frac{\log_e n}{n}}{p_n} > \frac{s}{2}$.

Proof: In the supplementary material. ■

Remark 4: Under certain sparsity conditions on graph edges, sufficient conditions for successful matching were derived in [20]. The erasure Model described above satisfies these sparsity conditions, and [20] provides guarantees for successful matching when $\lim_{n \rightarrow \infty} \frac{\frac{\log_e n}{n}}{p_n} \leq s$.

VI. MATCHING GRAPHS WITH COMMUNITY STRUCTURE

In this section, we describe the TM scheme for matching graphs generated under the SBM, i.e., graphs with community structure and provide achievable regions for these matching scenarios. A pair of correlated graphs with community structure are a special instance of the correlated random graphs defined in Definition 2. In order to describe the notation used in this section, we provide a separate formal definition of random graphs with community structure below.

A. Problem Setup

To describe the notation used in the section, consider a graph with $n \in \mathbb{N}$ vertices belonging to $c \in \mathbb{N}$ communities whose edges take $l \geq 2$ possible attributes. It is assumed that the set of communities $\mathcal{C} = \{C_1, C_2, \dots, C_c\}$ partitions the vertex set \mathcal{V} . The i^{th} community is written as $C_i = \{v_{j_1}, v_{j_2}, \dots, v_{j_{n_i}}\}$, where $n_i \in [n]$ is the size of the i^{th} community. Consequently, the graph is parametrized by $(n, c, (n_i)_{i \in [c]}, l)$. We sometimes refer to such an unlabeled graph as an $(n, c, (n_i)_{i \in [c]}, l)$ -unlabeled graph with community structure (UCS). The set $\mathcal{E}_{i_1, i_2} = \{(x, v_{j_1}, v_{j_2}) \in \mathcal{E} | v_{j_1} \in C_{i_1}, v_{j_2} \in C_{i_2}\}$ is the set of edges connecting the vertices in communities C_{i_1} and C_{i_2} . It can be noted that The Erdős-Rényi (ER) graphs studied in Section V are examples of single-community graphs, i.e., $c = 1$.

We consider random graphs with community structure (RCS) generated stochastically based on the SBM model. In this model, the probability of an edge between a pair of vertices is determined by their community memberships. More precisely, for a given vertex set \mathcal{V} and set of communities \mathcal{C} , let $P_{X|C_{j_1}, C_{j_2}, j_1, j_2 \in [c]}$ be a set of conditional distributions defined on $\mathcal{X} \times \mathcal{C} \times \mathcal{C}$, where $\mathcal{X} = [0, l - 1]$. It is assumed that the edge set \mathcal{E} is generated randomly, where the attribute X of the edge between vertices $v_{i_1} \in C_{j_1}$ and $v_{i_2} \in C_{j_2}$ is generated based on the conditional distribution $P_{X|C_{j_1}, C_{j_2}}$. So,

$$P((x, v_{i_1}, v_{i_2}) \in \mathcal{E}) = P_{X|C_{j_1}, C_{j_2}}(x | C_{j_1}, C_{j_2}), \quad \forall x \in [0, l - 1],$$

where $v_{i_1}, v_{i_2} \in C_{j_1} \times C_{j_2}$, and edges between different vertices are mutually independent. It can be noted that for undirected graphs considered in this work, we must have $P_{X|C_{j_1}, C_{j_2}}(x | C_{j_1}, C_{j_2}) = P_{X|C_{j_2}, C_{j_1}}(x | C_{j_2}, C_{j_1})$. A labeled graph with community structure is a graph with community structure g equipped with a labeling σ , and is denoted by $\tilde{g} = (g, \sigma)$. For the labeled graph \tilde{g} the adjacency matrix is defined as $G_\sigma = [G_{\sigma, i, j}]_{i, j \in [1, n]}$ where $G_{\sigma, i, j}$ is the unique value such that $(G_{\sigma, i, j}, v_k, v_l) \in \mathcal{E}_n$, where $(v_k, v_l) = (\sigma^{-1}(i), \sigma^{-1}(j))$. The submatrix $G_{\sigma, C_i, C_j} = [G_{\sigma, k, l}]_{k, l: v_k, v_l \in C_i \times C_j}$ is the adjacency matrix corresponding to the pair C_i and C_j . The upper triangle (UT) corresponding to \tilde{g} is the structure $U_\sigma = [G_{\sigma, i, j}]_{i < j}$. The upper triangle corresponding to communities C_i and C_j in \tilde{g} is denoted by $U_{\sigma, C_i, C_j} = [G_{\sigma, k, l}]_{k < l: v_k, v_l \in C_i \times C_j}$. The subscript

TABLE II
NOTATION TABLE: GRAPHS WITH COMMUNITY STRUCTURE

n :	# of vertices	c :	# of communities	C :	set of communities
C_i :	i^{th} community	n_i :	size of i^{th} community	$\mathcal{E}_{i,j}$:	edges between C_i and C_j
G_σ :	adjacency matrix	U_σ :	upper-triangle	$G_{\sigma,i,j}$:	adj. matrix between C_i and C_j

' σ ' is dropped when there is no ambiguity. The notation is summarized in Table II.

We consider pairs of correlated RCSs. It is assumed that edges between pairs of vertices in the two graphs with the same labeling are correlated and are generated based on a joint probability distribution, whereas edges between pairs of vertices with different labeling are generated independently. A pair of correlated RCSs is formally defined below.

Definition 18 (Correlated Pair of RCSs): Let $P_{X,X'|C_{j_1},C_{j_2},C'_{j_1},C'_{j_2}}, j_1,j_2,j'_1,j'_2 \in [1,c]$ be a set of conditional distributions defined on $\mathcal{X} \times \mathcal{X}' \times \mathcal{C} \times \mathcal{C}' \times \mathcal{C}'$, where $\mathcal{X} = \mathcal{X}' = [0, l-1]$ and $(\mathcal{C}, \mathcal{C}')$ are a pair of community sets of size $c \in \mathbb{N}$. A correlated pair of random graphs with community structure (CPCS) generated according to $P_{X,X'|C_{j_1},C_{j_2},C'_{j_1},C'_{j_2}}$ is a pair $\tilde{g} = (\tilde{g}, \tilde{g}')$ characterized by:

i) the pair of RCSs (g, g') generated according to $P_{X|C_{j_1},C_{j_2}}$ and $P_{X'|C'_{j_1},C'_{j_2}}$, respectively, ii) the pair of labelings (σ, σ') , and iii) the probability distribution $P_{X,X'|C_{j_1},C_{j_2},C'_{j_1},C'_{j_2}}$, such that:

- 1) The graphs have the same set of vertices $\mathcal{V} = \mathcal{V}'$.
- 2) For any two edges $e = (x, v_{j_1}, v_{j_2}), e' = (x', v'_{j_1}, v'_{j_2}), x, x' \in [0, l-1]$, we have
$$Pr(e \in \mathcal{E}; e' \in \mathcal{E}') = \begin{cases} P_{X,X'}(x, x'), & \text{if } \sigma(v_{j_k}) = \sigma'(v'_{j_k}), k = 1, 2 \\ Q_{X,X'}(x, x'), & \text{Otherwise,} \end{cases}$$

where $l \in \{1, 2\}$, $v_{j_1}, v_{j_2} \in C_{j_1} \times C_{j_2}, v'_{j_1}, v'_{j_2} \in C'_{j_1} \times C'_{j_2}$, the distribution $P_{X,X'}$ is the joint edge distribution when the edges connect vertices with similar labels and is given by $P_{X,X'|C_{j_1},C_{j_2},C'_{j_1},C'_{j_2}}$, the distribution $Q_{X,X'}$ is the conditional edge distribution when the edges connect labels with different labels and is given by $P_{X|C_{j_1},C_{j_2}} \times P_{X'|C'_{j_1},C'_{j_2}}$.

In this article, in order to simplify the notation, we assume that the community memberships in both graphs are the same. In other words, we assume that $v_j \in C_i \Rightarrow v'_{j'} \in C_i$ given that $\sigma(v_j) = \sigma'(v'_{j'})$ for any $j, j' \in [n]$ and $i \in [c]$. Furthermore, we assume that the size of the communities in the graph sequence grows linearly in the number of vertices. More precisely, let $\Lambda^{(n)}(i) \triangleq |C_i^{(n)}|$ be the size of the i^{th} community, we assume that³ $\Lambda^{(n)}(i) = \Theta(n)$ for all $i \in [c]$. We also assume that the number of communities c is constant in n .

We consider matching strategies under two scenarios:

- **With Complete Side-Information:** In this scenario, the matching strategy uses prior knowledge of vertices' community memberships. A matching strategy operating with complete side-information is a sequence of functions

³We write $f(x) = \Theta(g(x))$ if $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)}$ is a non-zero constant.

$f_n^{CSI} : (g^{(n)}, C^{(n)}, C'^{(n)}) \mapsto \hat{\sigma}'^{(n)}, n \in \mathbb{N}$, where $\underline{g}^{(n)} = (\underline{g}_1^{(n)}, \underline{g}_2^{(n)})$ consists of a pair of graphs with community structure with n vertices.

- **With Partial Side-Information:** A matching strategy operating with partial side-information does not use prior knowledge of the vertices' community memberships, rather, it uses the statistics $P_{X,X'|C_i,C_o,C'_{i'},C'_{o'}}$ and the community sizes $(n_i)_{i \in [c]}$. The matching strategy is a sequence of functions $f_n^{WSI} : \underline{g}^{(n)} \mapsto \hat{\sigma}'^{(n)}, n \in \mathbb{N}$.

The matching strategy is said to be successful if the fraction of correctly matched vertices approaches 1 as $n \rightarrow \infty$ as formalized in Definition 3.

B. Matching in Presence of Side-Information

First, we describe the matching strategy under the complete side-information scenario. In this scenario, the community membership of the nodes at both graphs are known prior to matching. Given a CPCS \tilde{g} generated according to $P_{X,X'|C_{j_1},C_{j_2},C'_{j_1},C'_{j_2}}, j_1,j_2,j'_1,j'_2 \in [1,c]$, the scheme operates as follows. It finds a labeling $\hat{\sigma}'$, for which i) the set of pairs $(G_{\sigma,C_{j_1},C_{j_2}}, G'_{\hat{\sigma}',C'_{j_1},C'_{j_2}}), j_1,j_2 \in [c]$ are jointly typical each with respect to $P_{X,X'|C_{j_1},C_{j_2},C'_{j_1},C'_{j_2}}(\cdot, \cdot | C_{j_1}, C_{j_2}, C'_{j_1}, C'_{j_2})$ when viewed as vectors of length $n_i n_j, i \neq j$, and ii) the set of pairs $(U_{\sigma,C_j,C_j}, U'_{\hat{\sigma}',C'_j,C'_j}), j \in [c]$ are jointly typical with respect to $P_{X,X'|C_{j_1},C_{j_2},C'_{j_1},C'_{j_2}}(\cdot, \cdot | C_j, C_j, C'_j, C'_j)$ when viewed as vectors of length $\frac{n_i(n_i-1)}{2}, j \in [c]$. Specifically, it returns a randomly picked element $\hat{\sigma}'$ from the set:

$$\begin{aligned} \hat{\Sigma}_{C,C'} = & \left\{ \hat{\sigma}' | \left(U_{\sigma,C_j,C_j}, U'_{\hat{\sigma}',C'_j,C'_j} \right) \in \mathcal{A}_\epsilon^{\frac{n_j(n_j-1)}{2}} \right. \\ & \times \left(P_{X,X'|C_j,C_j,C'_j,C'_j} \right), \forall j \in [c], \left(G_{\sigma,C_i,C_j}, G'_{\hat{\sigma}',C'_i,C'_j} \right) \\ & \left. \in \mathcal{A}_\epsilon^{n_i n_j} \left(P_{X,X'|C_i,C_j,C'_i,C'_j} \right), \forall i, j \in [c], i \neq j \right\}, \end{aligned}$$

where $\epsilon = \omega(\frac{1}{n})$, and declares $\hat{\sigma}'$ as the correct labeling. We show that under this scheme, the probability of incorrect labeling for any given vertex is arbitrarily small for large n .

Theorem 6: For the TM strategy described above, a given family of sets of distributions $\tilde{\mathcal{P}} = (\mathcal{P}^{(n)})_{n \in \mathbb{N}}$ is achievable, if for any constants $\delta > 0, \alpha \in [0, 1 - \delta]$ and every sequence of distributions $P_{X,X'|C_{j_1},C_{j_2},C'_{j_1},C'_{j_2}}^{(n)} \in \mathcal{P}_n, j_1,j_2,j'_1,j'_2 \in [1,c]$, and community sizes (n_1, n_2, \dots, n_c) such that $\sum_{i=1}^c n_i = n$, the

following holds:

$$3(1-\alpha)\frac{\log n}{n} \leq \min_{[\alpha_i]_{i \in [c]} \in \mathcal{A}_\alpha} \sum_{i,j \in [c], i < j} \frac{n_i n_j}{n^2} \\ \times D\left(P_{X,X'|C_i,C_j}^{(n)} \parallel (1-\beta_{i,j})P_{X|C_i,C_j}^{(n)} P_{X'|C_i,C_j}^{(n)} + \beta_{i,j}P_{X,X'|C_i,C_j}^{(n)}\right) \\ + \sum_{i \in [c]} \frac{n_i(n_i-1)}{2n^2} D\left(P_{X,X'|C_i,C_i}^{(n)} \parallel (1-\beta_i)P_{X|C_i,C_i}^{(n)} \right. \\ \left. \times P_{X'|C_i,C_i}^{(n)} + \beta_i P_{X,X'|C_i,C_i}^{(n)}\right), \quad (19)$$

and $\max_{(x_1,x_2): P_{X_1|C_i}(x_1)P_{X_2|C_j}(x_2) \neq 0} |\log \frac{P_{X_1|C_i}(x_1)P_{X_2|C_j}(x_2)}{P_{X_1,X_2|C_i,C_j}(x_1,x_2)}| + = o(\log n)$, $i, j \in [c]$, as $n \rightarrow \infty$, where $\mathcal{A}_\alpha = \{([\alpha_i]_{i \in [c]}): \alpha_i \leq \frac{n_i}{n}, \sum_{i \in [c]} \alpha_i = \alpha\}$, and $\beta_{i,j} = \frac{n_i^2}{n_i n_j} \alpha_i \alpha_j$, $i, j \in [c]$ and $\beta_i = \frac{n \alpha_i (n \alpha_i - 1)}{n_i (n_i - 1)}$, $i \in [c]$. The maximal family of sets of distributions which are achievable using the typicality matching strategy with complete side-information is denoted by \mathcal{P}_{full} .

Proof: In the supplementary material. ■

Remark 5: Note that the community sizes (n_1, n_2, \dots, n_c) , $n \in \mathbb{N}$ are assumed to grow in n such that $\lim_{n \rightarrow \infty} \frac{n_i}{n} > 0$.

It can be noted that Theorem 6 includes the achievable region for matching of pairs of Erdős-Rényi graphs (i.e., single community) derived in Theorem 3.

C. Matching in Absence of Side-Information

The scheme described in the previous section can be extended to matching graphs without community memberships side-information. In this scenario, it is assumed that the distribution $P_{X,X'|C_{j_1}, C_{j_2}, C'_{j_1}, C'_{j_2}}, j_1, j_2, j'_1, j'_2 \in [1, c]$ is known, but the community memberships of the vertices in the graphs are not known. In this case, the scheme sweeps over all possible community membership assignments of the vertices in the two graphs. For each community membership assignment, the scheme attempts to match the two graphs using the method proposed in the complete side-information scenario. If it finds a labeling which satisfies the joint typicality conditions, it declares the labeling as the correct labeling. Otherwise, the scheme proceeds to the next community membership assignment. More precisely, for a given community assignment (\hat{C}, \hat{C}') , the scheme forms the following ambiguity set

$$\hat{\Sigma}_{\hat{C}, \hat{C}'} = \left\{ \hat{\sigma}' \mid \left(U_{\sigma, \hat{C}_i, \hat{C}_{i4}}, U'_{\hat{\sigma}', \hat{C}'_i, \hat{C}'_{i4}} \right) \in \mathcal{A}_\epsilon^{\frac{n_i(n_i-1)}{2}} \left(P_{X,X'|C_i, C_i, C'_i, C'_i} \right), \right. \\ \left. \forall i \in [c], \left(G_{\sigma, \hat{C}_i, \hat{C}_j}, \tilde{G}'_{\hat{\sigma}', \hat{C}'_i, \hat{C}'_j} \right) \right. \\ \left. \in \mathcal{A}_\epsilon^{n_i n_j} \left(P_{X,X'|C_i, C_j, C'_i, C'_j} \right), \quad \forall i, j \in [c], i \neq j \right\}.$$

Define $\hat{\Sigma}_0 \triangleq \bigcup_{(\hat{C}, \hat{C}') \in \mathcal{C}} \hat{\Sigma}_{\hat{C}, \hat{C}'}$, where \mathcal{C} is the set of all possible community membership assignments. The scheme outputs a randomly chosen element of $\hat{\Sigma}_0$ as the correct labeling. The following theorem shows that the achievable region is the same as the one described in Theorem 6.

Theorem 7: Let \mathcal{P}_0 be the maximal family of sets of achievable distributions for the typicality matching strategy without side-information. Then, $\mathcal{P}_0 = \mathcal{P}_{full}$.

The proof follows similar arguments as Theorem 6. We provide an outline. It is enough to show that $|\hat{\Sigma}_0|$ has the same exponent as $|\hat{\Sigma}_{C,C'}|$. Note that the size of the set of all community membership assignments \mathcal{C} has an exponent which is $\Theta(n)$ since $|\mathcal{C}| \leq 2^{cn}$. On the other hand,

$$|\hat{\Sigma}_0| \leq |\mathcal{C}| |\hat{\Sigma}_{C,C'}| \leq 2^{nc} 2^{\Theta(n \log n)} = 2^{\Theta(n \log n)}.$$

The rest of the proof follows by the same arguments as in Theorem 6.

VII. MATCHING COLLECTIONS OF GRAPHS

In the previous sections, we considered matching of pairs of correlated graphs. The results can be further extended to problems involving matching of collections of more than two graphs. In this section, we consider matching collections of more than two correlated graphs, where the first graph is deanonymized and the other graphs are anonymized. For brevity we consider collections of correlated Erdős-Rényi graphs, i.e., single-community random graphs in Section V. The results can be further extended to correlated graphs with community structure in a straightforward manner. We formally describe collections of correlated Erdős-Rényi graphs below.

Definition 19 (Correlated Collection of ER Graphs): Let P_{X^m} be a conditional distribution defined on $\prod_{k \in [m]} \mathcal{X}_k$, where $\mathcal{X}_i = [0, l-1]$, $i \in [m]$ and $m > 2$. A correlated collection of ER graphs $\tilde{g} = (\tilde{g}^i)_{i \in [m]}$ generated according to P_{X^m} is characterized by: i) the collection of ER graphs $(g^i)_{i \in [m]}$ each generated according to P_{X_i} , ii) the collection of labelings $(\sigma_i)_{i \in [m]}$ for the unlabeled graphs $(g^i)_{i \in [m]}$, and iii) the joint probability distribution P_{X^m} , such that:

- 1) The graphs have the same set of vertices $\mathcal{V} = \mathcal{V}_i$, $i \in [m]$.
- 2) For any collection of edges $e^i = (x^i, v_{j_1}^i, v_{j_2}^i)$, $x^i \in [0, l-1]$, $i \in [m]$, we have

$$Pr(e^i \in \mathcal{E}^i, i \in [m]) \\ = \begin{cases} P_{X^m}(x^m), & \text{if } \sigma^i(v_{j_1}^i) = \sigma^k(v_{j_1}^k), \quad \forall i, k \in [m] \\ \prod_{i \in [m]} P_{X_i}(x_i), & \text{Otherwise,} \end{cases}$$

where $l \in \{1, 2\}$, and $v_{j_1}^i, v_{j_2}^i \in \mathcal{V}_1 \times \mathcal{V}_2$, $i \in [m]$.

Similar to the TM strategy for pairs of correlated graphs described in Section V, we propose a matching strategy based on typicality for collections of correlated graphs. Given a correlated collection of graphs $(g^i)_{i \in [m]}$, where the labeling for \tilde{g}^1 is given and the rest of the graphs are anonymized, the TM strategy operates as follows. The scheme finds a collection $\hat{\Sigma}$ of labelings $\hat{\sigma}^j$, $j \in [2, m]$, for which the UT's $U_{\sigma^j}^j$, $j \in [m]$ are jointly typical with respect to P_{n, X^m} when viewed as vectors of length $\frac{n(n-1)}{2}$. The strategy succeeds if at least one such labeling exists and fails otherwise.

Theorem 8: For the TM strategy, a given family of sets of distributions $\tilde{P} = (\mathcal{P}_n)_{n \in \mathbb{N}}$ is achievable, if for every sequence

of distributions $P_{n,X^m} \in \mathcal{P}_n, n \in \mathbb{N}$ we have

$$\frac{\log n}{n} \left(\sum_{k \in [b_m]} |\mathcal{P}_k| \alpha_k - 1 \right) \leq \frac{1}{2(b_m - 1)(m(m - 1) + 1)} \\ \times D \left(P_{X^m} \parallel \sum_{k \in [b_m]} \alpha'_k P_{X_{\mathcal{P}_k}} \right) + O \left(\frac{\log n}{n} \right), \quad (20)$$

for all $\alpha_1, \alpha_2, \dots, \alpha_{b_m} : \sum_{k \in [b_m]} \alpha_k = n, \alpha_{b_m} \in [1, 1 - \alpha_n]$, and $\max_{X^m: P_{X^m}(x^m) \neq 0} |\log \frac{\prod_{i \in [m]} P_{X_i}(x_i)}{P_{X^m}(x^m)}| = o(\log n)$, where $\alpha'_k = \frac{\alpha_k^2}{2} + \sum_{k', k'' : \mathcal{P}_{k'} \cap \mathcal{P}_{k''} = \mathcal{P}_k} \alpha_{k'} \alpha_{k''}$, $\mathcal{P}_{k'} = \{\mathcal{A}' \cap \mathcal{A}'' : \mathcal{A}' \in \mathcal{P}_{k'}, \mathcal{A}'' \in \mathcal{P}_{k''}\}$, $k', k'' \in [b_m]$, and $\mathcal{P}_{b_m} = [1, n]$.

Proof: In the supplementary material. ■

Remark 6: Note that Equation (20) recovers the result given in Equation (17) for matching of pairs of correlated ER graphs, i.e., $m = 2$.

VIII. CONVERSE RESULTS

In this section, we provide conditions on the graph parameters under which graph matching is not possible. Without loss of generality, we assume that (σ, σ') are a pair of random labelings chosen uniformly among the set of all possible labelings for the two graphs. Roughly speaking, the information revealed by identifying the realization of σ' is equal to $H(\sigma') = \log(n!) \approx \log(n^n) = n \log n$. Consequently, using Fano's inequality, we show that the information contained in (σ, g, g') regarding σ' , which is quantified as the mutual information $I(\sigma'; \sigma, g, g')$, must be at least $n \log n$ bits for successful matching. The mutual information $I(\sigma'; \sigma, g, g')$ is a function of multi-letter probability distributions. We use standard information theoretic techniques to bound $I(\sigma'; \sigma, g, g')$ using information quantities which are functionals of single-letter distributions. The following states the resulting necessary conditions for successful matching.

Theorem 9: For the graph matching problem under the community structure model with complete side-information, the following provides necessary conditions for successful matching:

$$\frac{\log n}{n} \leq \sum_{i,j \in [c], i < j} \frac{n_i n_j}{n^2} I(X, X' | C_i, C_j, C'_i, C'_j) \\ + \sum_{i \in [c]} \frac{n_i(n_i - 1)}{2n^2} I(X, X' | C_i, C_i, C'_i, C'_i) + o \left(\frac{\log n}{n} \right),$$

where $I(X, X' | C_i, C_j, C'_i, C'_j)$ is defined with respect to $P_{X, X' | C_i, C_j, C'_i, C'_j}$.

Proof: In the supplementary material. ■

Corollary 4: For the graph matching problem under the Erdős-Rényi model, the following provides necessary conditions for successful matching:

$$\frac{2 \log n}{n} \leq I(X, X') + o \left(\frac{\log n}{n} \right).$$

IX. SEEDED GRAPH MATCHING

So far, we have investigated the fundamental limits of graph matching assuming the availability of unlimited computational

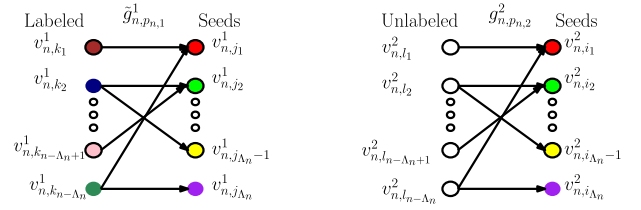


Fig. 3. The matching algorithm constructs the bipartite graph which captures the connections between the unmatched vertices with the seed vertices.

resources. In this section, we consider seeded graph matching, and propose a matching algorithm whose complexity grows polynomially in the number of vertices of the graph and leads to successful matching in a wide range of graph matching scenarios. The algorithm leverages ideas from prior work a related problem called *online fingerprinting* which involves matching of correlated bipartite graphs [46].

In seeded graph matching, it is assumed that we are given the correct labeling for a subset of the vertices in the anonymized graph prior to the start of the matching process. The subset of pre-matched vertices are called ‘seeds’. The motivation behind the problem formulation is that in many applications of graph matching, the correct labeling of a subset of vertices is known through side-information. For instance, in social network deanonymization, many users link their social media accounts across social networks publicly. As shown in this section, the seed side-information can be used to significantly reduce the complexity of the matching algorithm.

The proposed graph matching algorithm operates as follows. First, the algorithm constructs the bipartite graph shown in Figure 3 whose edges consist of the connections between the unmatched vertices with the seeded vertices in each graph. The algorithm proceeds in two steps. First, it constructs the ‘fingerprint’ vectors for each of the unmatched vertices in the two bipartite graphs based on their connections to the seed vertices. The fingerprint vector of a vertex is the row in the adjacency matrix of the bipartite graph corresponding to the edges between that vertex and the seed vertices. In the second step, the algorithm finds a jointly typical pair of fingerprint vectors in the deanonymized and deanonymized graph adjacency matrices and matches the corresponding vertices, where typicality is defined based on the joint distribution between the edges of the two graphs. Note that the bipartite graphs encompass only a subset of the edges in the original graphs. Hence by restricting the matching process to the bipartite graphs, some of the information which could potentially help in matching is ignored. This leads to more restrictive conditions on successful matching compared to the ones derived in the previous sections. However, the computational complexity of the resulting matching algorithm is considerably improved. In the following, we focus on matching of seeded CPERs. The results can be easily extended to seeded CPCs similar to the unseeded graph matching in prior sections. A seeded CPER (SCPER) is formally defined below.

Definition 20 (Correlated Pair of Seeded ER Graphs): An SCPER is a triple $(\tilde{g}, \tilde{g}', \mathcal{S})$, where $\tilde{g} = (\tilde{g}, \tilde{g}')$ is a CPER generated according to $P_{X, X'}$, and $\mathcal{S} \subseteq \mathcal{V}$ is the seed set.

Let $\mathcal{S} = \{v_{i_1}, v_{i_2}, \dots, v_{i_\Lambda}\}$ and define the reverse seed set $\mathcal{S}^{-1} = \{v_{j_1}, v_{j_2}, \dots, v_{j_\Lambda}\}$, where $\sigma(v_{j_k}) = \sigma'(v_{i_k}), k \in [1, \Lambda]$. The algorithm is given the correct labeling of all the vertices in the first graph $\sigma : \mathcal{V} \rightarrow [1, n]$ and the seed vertices in the second graph $\sigma'|_{\mathcal{S}} : \mathcal{S} \rightarrow [1, n]$. The objective is to find the correct labeling of the rest of the vertices in the second graph $\hat{\sigma}_n : \mathcal{V} \rightarrow [1, n]$ so that the fraction of mislabeled vertices is negligible as the number of vertices grows asymptotically large, i.e., $P(\hat{\sigma}' = \sigma') \rightarrow 1$ as $n \rightarrow \infty$. To this end, the algorithm first constructs a fingerprint for each vertex in each of the graphs. For an arbitrary vertex v_i in g_{P_X} , its fingerprint is defined as $\underline{F}_i = (F_i(1), F_i(2), \dots, F_i(\Lambda))$, which indicates its connections to the reverse seed elements:

$$F_i(l) = \begin{cases} 1 & \text{if } (v_i, v_{j_l}) \in \mathcal{E} \\ 0 & \text{Otherwise,} \end{cases} \quad l \in [1, \Lambda].$$

The fingerprint of a vertex v_i in the second graph is defined in a similar fashion based on connections to the elements of the seed set \mathcal{S} . Take an unmatched vertex $v_i \notin \mathcal{S}$. The algorithm matches v_i in g to a vertex v_j in g' if it is the unique vertex such that the fingerprint pair $(\underline{F}_i, \underline{F}'_j)$ are jointly ϵ -typical with respect to the distribution $P_{X, X'}$, where⁴ $\epsilon = \omega(\frac{1}{\sqrt{\Lambda}})$:

$$\exists! i : (\underline{F}_i, \underline{F}'_j) \in \mathcal{A}_\epsilon^n(X, X') \Rightarrow \hat{\sigma}(v_i) = \sigma'(v_j),$$

where $\mathcal{A}_\epsilon^n(X, X')$ is the set of jointly ϵ -typical set sequences of length n with respect to $P_{X, X'}$. If a unique match is not found, then vertex v_i is added to the ambiguity set \mathcal{L} . Hence, $\mathcal{V} \setminus \mathcal{L}$ is the set of all matched vertices. In the next step, these vertices are added to the seed set and the expanded seed set is used to match the vertices in the ambiguity set. The algorithm succeeds if all vertices are matched at this step and fails otherwise. We call this strategy the Seeded Typicality Matching Strategy (STM).

Theorem 10: Define the family of sets of pairs of distribution and seed sizes $\tilde{\mathcal{P}}$ as follows:

$$\tilde{\mathcal{P}} = \left\{ (\mathcal{P}_n, \Lambda_n)_{n \in \mathbb{N}} \mid \forall P_{n, X, X'} \in \mathcal{P}_n : \frac{2 \log n}{I(X, X')} \leq \Lambda_n, I(X; X') = \omega\left(\sqrt{\frac{1}{\Lambda_n}}\right) \right\}.$$

Any family of SCPEs with parameters chosen from $\tilde{\mathcal{P}}$ is matchable using the STM strategy.

The proof which is provided in the supplementary material uses the following lemma.

Lemma 6: The following holds:

$$P\left(|\mathcal{L}| > \frac{2n}{\Lambda \epsilon^2}\right) \rightarrow 0, \text{ as } n \rightarrow \infty,$$

Proof: In the supplementary material. ■

X. CONCLUSION

We have considered matching of collections of correlated graphs. We have studied the problem under the Erdős-Rényi

model as well as the more general community structure model. The derivations apply to graphs whose edges may take non-binary attributes. We have introduced a graph matching scheme called the Typicality Matching scheme which relies on tools such as concentration of measure and typicality of sequences of random variables to perform graph matching. We have further provided converse results which lead to necessary conditions on graph parameters for successful matching. We have investigated seeded graph matching, where the correct labeling of a subset of graph vertices is known prior to the matching process. We have introduced a matching algorithm for seeded graph matching which successfully matches the graphs in wide range of matching problems with large enough seeds and has a computational complexity which grows polynomially in the number of graph vertices.

REFERENCES

- [1] R. Gross and A. Acquisti, "Information revelation and privacy in online social networks," in *Proc. ACM Workshop Privacy Electron. Soc.*, 2005, pp. 71–80.
- [2] C. Shah, "Collaborative information seeking," *J. Assoc. Inf. Sci. Technol.*, vol. 65, no. 2, pp. 215–236, 2014.
- [3] L. Backstrom, C. Dwork, and J. Kleinberg, "Wherefore art thou r3579x?: Anonymized social networks, hidden patterns, and structural steganography," in *Proc. 16th Int. Conf. World Wide Web*, 2007, pp. 181–190.
- [4] A. Narayanan and V. Shmatikov, "De-anonymizing social networks," in *Proc. 30th IEEE Symp. Security Privacy*, Oakland, CA, USA, 2009, pp. 173–187.
- [5] S. Ji, W. Li, N. Z. Gong, P. Mittal, and R. Beyah, "Seed-based de-anonymizability quantification of social networks," *IEEE Trans. Inf. Forensics Security*, vol. 11, pp. 1398–1411, 2016.
- [6] E. Kazemi, L. Yartseva, and M. Grossglauser, "When can two unlabeled networks be aligned under partial overlap?" in *Proc. 53th Annu. Allerton Conf. Commun. Control Comput. (Allerton)*, Monticello, IL, USA, Sep. 2015, pp. 33–42.
- [7] F. Shirani, S. Garg, and E. Erkip, "Typicality matching for pairs of correlated graphs," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Vail, CO, USA, 2018, pp. 221–225.
- [8] P. Foggia, G. Percannella, and M. Vento, "Graph matching and learning in pattern recognition in the last 10 years," *Int. J. Pattern Recognit. Artif. Intell.*, vol. 28, no. 01, 2014, Art. no. 1450001.
- [9] K. Xu *et al.*, "Cross-lingual knowledge graph alignment via graph matching neural network," 2019. [Online]. Available: arXiv:1905.11605.
- [10] E. Kazemi, H. Hassani, M. Grossglauser, and H. P. Modarres, "PROPER: Global protein interaction network alignment through percolation matching," *BMC Bioinform.*, vol. 17, no. 1, p. 527, 2016.
- [11] P. Erdos and A. Rényi, "On the evolution of random graphs," *Publ. Math. Inst. Hung. Acad. Sci.*, vol. 5, no. 1, pp. 17–60, 1960.
- [12] E. M. Wright, "Graphs on unlabelled nodes with a given number of edges," *Acta Mathematica*, vol. 126, no. 1, pp. 1–9, 1971.
- [13] L. Babai, P. Erdos, and S. M. Selkow, "Random graph isomorphism," *SIAM J. Comput.*, vol. 9, no. 3, pp. 628–635, 1980.
- [14] B. Bollobás, *Random Graphs* (Cambridge Studies in Advanced Mathematics). Cambridge, U.K.: Cambridge Univ. Press, 2001.
- [15] T. Czajka and G. Pandurangan, "Improved random graph isomorphism," *J. Discrete Algorithms*, vol. 6, no. 1, pp. 85–92, 2008.
- [16] E. Kazemi, "Network alignment: Theory, algorithms, and applications," Ph.D. dissertation, Dept. Lab. Comput. Commun. Appl., EPFL, Lausanne, Switzerland, 2016.
- [17] L. Yartseva and M. Grossglauser, "On the performance of percolation graph matching," in *Proc. 1st ACM Conf. Online Soc. Netw.*, 2013, pp. 119–130.
- [18] P. Pedarsani, D. R. Figueiredo, and M. Grossglauser, "A Bayesian method for matching two similar graphs without seeds," in *Proc. IEEE 51st Annu. Allerton Conf. Commun. Control Comput. (Allerton)*, Monticello, IL, USA, 2013, pp. 1598–1607.
- [19] S. Ji, W. Li, M. Srivatsa, and R. Beyah, "Structural data de-anonymization: Quantification, practice, and implications," in *Proc. ACM SIGSAC Conf. Comput. Commun. Security*, 2014, pp. 1040–1053.

⁴Alternatively, $\lim_{n \rightarrow \infty} \frac{\epsilon}{\sqrt{|\mathcal{S}|}} = 0$.

- [20] D. Cullina and N. Kiyavash, "Exact alignment recovery for correlated Erdős-Rényi graphs," 2017. [Online]. Available: arXiv:1711.06783.
- [21] V. Lyzinski, "Information recovery in shuffled graphs via graph matching," 2016. [Online]. Available: arXiv:1605.02315.
- [22] D. Cullina, N. Kiyavash, P. Mittal, and H. V. Poor, "Partial recovery of Erdős-Rényi graph alignment via k -core alignment," 2018. [Online]. Available: arXiv:1809.03553.
- [23] M. Girvan and M. E. J. Newman, "Community structure in social and biological networks," *Proc. Nat. Acad. Sci.*, vol. 99, no. 12, pp. 7821–7826, 2002.
- [24] S. Fortunato and C. Castellano, "Community structure in graphs," 2007. [Online]. Available: arXiv:0712.2716.
- [25] F. Shirani, S. Garg, and E. Erkip, "Matching graphs with community structure: A concentration of measure approach," in *Proc. IEEE 56th Annu. Allerton Conf. Commun. Control Comput. (Allerton)*, Monticello, IL, USA, 2018, pp. 1028–1035.
- [26] S. Nilizadeh, A. Kapadia, and Y.-Y. Ahn, "Community-enhanced de-anonymization of online social networks," in *Proc. ACM SIGSAC Conf. Comput. Commun. Security*, 2014, pp. 537–548.
- [27] K. Singhal, D. Cullina, and N. Kiyavash, "Significance of side information in the graph matching problem," 2017. [Online]. Available: arXiv:1706.06936.
- [28] V. Lyzinski, D. L. Sussman, D. E. Fishkind, H. Pao, J. T. Vogelstein, and C. E. Priebe, "Seeded graph matching for large stochastic block model graphs," *Stat.*, vol. 1050, p. 12, 2014.
- [29] E. Onaran, S. Garg, and E. Erkip, "Optimal de-anonymization in random graphs with community structure," in *Proc. 50th Asilomar Conf. Signals Syst. Comput.*, Newark, NJ, USA, 2016, pp. 709–713.
- [30] D. Cullina and N. Kiyavash, "Improved achievability and converse bounds for erdos-renyi graph matching," *SIGMETRICS Perform. Eval. Rev.*, vol. 44, no. 1, pp. 63–72, Jun. 2016.
- [31] E. Kazemi, S. H. Hassani, and M. Grossglauser, "Growing a graph matching from a handful of seeds," *Proc. VLDB Endowm.*, vol. 8, no. 10, pp. 1010–1021, 2015.
- [32] C.-F. Chiasserini, M. Garetto, and E. Leonardi, "Social network de-anonymization under scale-free user relations," *IEEE/ACM Trans. Netw.*, vol. 24, no. 6, pp. 3756–3769, Dec. 2016.
- [33] V. Lyzinski, D. E. Fishkind, and C. E. Priebe, "Seeded graph matching for correlated Erdős-Rényi graphs," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 3513–3540, 2014.
- [34] M. Fiori, P. Sprechmann, J. Vogelstein, P. Musé, and G. Sapiro, "Robust multimodal graph matching: Sparse coding meets graph matching," in *Advances in Neural Information Processing Systems*. Red Hook, NY, USA: Curran Assoc., Inc., 2013, pp. 127–135.
- [35] F. Shirani, S. Garg, and E. Erkip, "Seeded graph matching: Efficient algorithms and theoretical guarantees," in *Proc. IEEE 51st Asilomar Conf. Signals Syst. Comput.*, Pacific Grove, CA, USA, 2017, pp. 253–257.
- [36] E. Mossel and J. Xu, "Seeded graph matching via large neighborhood statistics," in *Proc. 13th Annu. ACM-SIAM Symp. Discrete Algorithms*, 2019, pp. 1005–1014.
- [37] D. E. Fishkind *et al.*, "Seeded graph matching," *Pattern Recognit.*, vol. 87, pp. 203–215, Mar. 2019.
- [38] V. Lyzinski and D. L. Sussman, "Matchability of heterogeneous networks pairs," 2017. [Online]. Available: arXiv:1705.02294.
- [39] S. Zhang and H. Tong, "FINAL: Fast attributed network alignment," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Min.*, 2016, pp. 1345–1354.
- [40] M. Heimann, H. Shen, T. Safavi, and D. Koutra, "REGAL: Representation learning-based graph alignment," in *Proc. 27th ACM Int. Conf. Inf. Knowl. Manag.*, 2018, pp. 117–126.
- [41] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*. New York, NY, USA: Academic, 1981.
- [42] I. M. Isaacs, *Algebra: A Graduate Course*, vol. 100. Providence, RI, USA: Amer. Math. Soc., 1994.
- [43] E. Tuncel, "On error exponents in hypothesis testing," *IEEE Trans. Inf. Theory*, vol. 51, no. 8, pp. 2945–2950, Aug. 2005.
- [44] I. Csiszár, "The method of types [information theory]," *IEEE Trans. Inf. Theory*, vol. 44, no. 6, pp. 2505–2523, Oct. 1998.
- [45] L. Comtet, *Advanced Combinatorics: The Art of Finite and Infinite Expansions*. Dordrecht, The Netherlands: Springer, 2012.
- [46] F. Shirani, S. Garg, and E. Erkip, "An information theoretic framework for active de-anonymization in social networks based on group memberships," in *Proc. 55rd Annu. Allerton Conf. Commun. Control Comput.*, Sep. 2017, pp. 470–477.



Farhad Shirani (Member, IEEE) received the B.S. degree in electrical engineering from the Sharif University of Technology, the M.Sc. degree in applied mathematics and the Ph.D. degree from the University of Michigan at Ann Arbor, Ann Arbor, MI, USA. He is an Assistant Professor with the Electrical and Computer Engineering Department, North Dakota State University. He served as a Lecturer and a Postdoctoral Research Fellow with the University of Michigan in 2017. He was a Research Assistant Professor with New York University from 2017 to 2020. His research interests include privacy and security, wireless communications, and machine learning. His recent work include developing information theoretic methods for analysis of fundamental limits of Web privacy, design of receiver architectures for energy efficient communication over MIMO systems, and design of algorithms for opportunistic multi-user scheduling under various fairness constraints.



Siddharth Garg received the B.Tech. degree in electrical engineering from the Indian Institute of Technology Madras, and the Ph.D. degree in electrical and computer engineering from Carnegie Mellon University in 2009.

He is currently an Institute Associate Professor of ECE with the Tandon School of Engineering, New York University, where he leads the EnSuRe Research. He was an Assistant Professor in ECE with the University of Waterloo from 2010 to 2020.

His research interests are in machine learning, cyber-security, and computer hardware design. In 2016, he was listed in Popular Science Magazine's annual list of "Brilliant 10" researchers. He has received the NSF CAREER Award in 2015, and Paper Awards at the IEEE Symposium on Security and Privacy in 2016, the USENIX Security Symposium in 2013, at the Semiconductor Research Consortium TECHCON in 2010, and the International Symposium on Quality in Electronic Design in 2009. He also received the Angel G. Jordan Award from ECE Department of Carnegie Mellon University for outstanding thesis contributions and service to the community. He serves on the technical program committee of several top conferences in the area of computer engineering and computer hardware, and has served as a reviewer for several IEEE and ACM journals.



Elza Erkip (Fellow, IEEE) received the B.S. degree in electrical and electronics engineering from Middle East Technical University, Ankara, Turkey, and the M.S. and Ph.D. degrees in electrical engineering from Stanford University, Stanford, CA, USA. She is an Institute Professor with the Electrical and Computer Engineering Department, Tandon School of Engineering, New York University. Her research interests are in information theory, communication theory, and wireless communications. She is a Member of the Science Academy of Turkey and is a

Clarivate Highly Cited Researcher. She received the NSF CAREER Award in 2001, the IEEE Communications Society WICE Outstanding Achievement Award in 2016, and the IEEE Communications Society Communication Theory Technical Committee Technical Achievement Award in 2018. Her paper awards include the IEEE Communications Society Stephen O. Rice Paper Prize in 2004, and the IEEE Communications Society Award for Advances in Communication in 2013. She has been a Member of the Board of Governors of the IEEE Information Theory Society since 2012, where in 2018 she was the Society President. She was a Distinguished Lecturer of the IEEE Information Theory Society from 2013 to 2014. She has had many editorial and conference organization responsibilities. Some recent ones include, the Asilomar Conference on Signals, Systems and Computers, MIMO Communications and Signal Processing Track Chair in 2017, the IEEE Wireless Communications and Networking Conference Technical Co-Chair in 2017, the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS Guest Editor in 2015, and the IEEE International Symposium of Information Theory General Co-Chair in 2013.