


RESEARCH

Open Access



The phantom alignment strength conjecture: practical use of graph matching alignment strength to indicate a meaningful graph match

Donniell E. Fishkind^{1*} , Felix Parker², Hamilton Sawczuk¹, Lingyao Meng¹, Eric Bridgeford³, Avanti Athreya¹, Carey Priebe¹ and Vince Lyzinski⁴

*Correspondence:
def@jhu.edu

¹ Department of Applied
Mathematics and Statistics,
Johns Hopkins University,
Baltimore, MD, USA
Full list of author information
is available at the end of the
article

Abstract

The alignment strength of a graph matching is a quantity that gives the practitioner a measure of the correlation of the two graphs, and it can also give the practitioner a sense for whether the graph matching algorithm found the true matching. Unfortunately, when a graph matching algorithm fails to find the truth because of weak signal, there may be “phantom alignment strength” from meaningless matchings that, by random noise, have fewer disagreements than average (sometimes substantially fewer); this alignment strength may give the misleading appearance of significance. A practitioner needs to know what level of alignment strength may be phantom alignment strength and what level indicates that the graph matching algorithm obtained the true matching and is a meaningful measure of the graph correlation. The *Phantom Alignment Strength Conjecture* introduced here provides a principled and practical means to approach this issue. We provide empirical evidence for the conjecture, and explore its consequences.

Keywords: Graph matching, Alignment strength, Graph correlation

Introduction

This paper is about graph matchability in practice. Specifically, when given two graphs and an unobserved “true” bijection (also called “true matching” or “true alignment”) between their vertices, will exact (i.e. optimal) graph matching and approximate graph matching algorithms provide us with the matching which is the “truth”? How might we know in actual practice whether the “truth” has been found? Our work is in response to the latter question. The main contribution here is our formulation of the *Phantom Alignment Strength Conjecture* in Section “[Phantom alignment strength conjecture, consequences](#)”, followed up in Sect. “[Phantom alignment strength conjecture, consequences](#)” with the practical implications of this conjecture in deciding when alignment strength is high enough to indicate truth. This conjecture is also interesting as a theoretical matter, completely aside from its consequences.

Graphs (networks) are a commonly used data modality for encoding relationships, interactions, and dependencies in data in an incredibly broad range of the sciences and engineering; this includes sociology (e.g., social network analysis Wasserman and Faust 1994), neuroscience connectomics (Bullmore and Sporns 2009; Vogelstein et al. 2019), biology (e.g., biological interaction networks Szklarczyk et al. 2015; Tong et al. 2004), and automated knowledge discovery (Wang et al. 2017), to name just a few application areas.

The graph matching problem is, given two graphs with the same number of vertices, to find the bijection between the vertex sets that minimizes the number of adjacency “disagreements” between the graphs. Often there is an underlying “true” bijection that the graph matching is attempting to recover/approximate. Sometimes part of this true bijection is known a-priori, in which case minimizing the number of disagreements over the remainder of the bijection is called seeded graph matching. Graph matching and seeded graph matching are formally defined in Sect. “[Overview: seeded graph matching, alignment strength](#)”.

Graph matching and seeded graph matching are used in a wide variety of places, and we mention just a few. Information about the interactions amongst objects of interest is sometimes split across multiple networks or multiple layers of the same network (Kivelä et al. 2014). In many applications, such as neuroscience connectomics where, for example, DT-MRI derived graphs can be generated by aligning scans to a common template before uncovering the underlying edge structure (Gray et al. 2012), the vertices across networks or across layers are a priori aligned and identified. These aligned vertex labels can then be used to create joint network inference procedures that can leverage the signal across multiple networks for more powerful statistical inference (Levin et al. 2017; Chen et al. 2016; Arroyo et al. 2019; Durante and Dunson 2018). In many other applications, the vertex labels across networks or across layers are unknown or noisily observed. Social networks provide a canonical example of this, where common users across different social network platforms may use different user names and their user profiles may not be linked across networks. Discovering this latent correspondence (in the social network example, this is anchoring profiles to a common user across networks) is a key inference task (Lin et al. 2010; Yartseva and Grossglauser 2013) for leveraging the information across networks for subsequent inference, and it is a key consideration for understanding the degree of user anonymity (Ding et al. 2010) across platforms.

For a thorough survey of the relevant graph matching literature, see Conte et al. (2004), Foggia et al. (2014), Emmert-Streib et al. (2016).

The graph matching problem is computationally complex. Indeed, the simpler graph isomorphism problem has been shown to be of quasi-polynomial complexity (Babai 2016). Allowing loopy, weighted, directed graphs makes graph matching equivalent to the NP-hard quadratic assignment problem. Due to its practical importance and computational difficulty, a large branch of the graph matching literature is devoted to developing algorithms to efficiently, but approximately, solve the graph matching problem; see, for example, Fishkind et al. (2019a), Umeyama (1988), Singh et al. (2007), Zaslavskiy et al. (2009), Zhou and De la Torre (2012), Vogelstein et al. (2014), Zhang and Tong (2016), Feizi et al. (2016), Heimann et al. (2018) among myriad others.

Somewhat dual to the algorithmic development literature, a large branch of the modern graph matching literature is devoted to theoretically exploring the question of graph matchability, also called graph de-anonymization; this is the question of determining when there is enough signal present for graph matching to recover the “true” bijection. Many of the recent papers in this area have introduced latent alignment across graphs by correlating the edges across networks between common pairs of vertices, focusing on understanding the phase transition between matchable and non-matchable networks in terms of the level of correlation across networks and/or the sparsity level of the networks; see, for example, Pedarsani and Grossglauser (2011), Patsolic et al. (2014), Cullina and Kiyavash (2016), Lyzinski et al. (2016), Cullina and Kiyavash (2017), Sussman et al. (2019), Cullina et al. (2019), Fan et al. (2019), Ding et al. (2020), Mossel and Xu (2020).

In Fishkind et al. (2019b), a novel measure of graph correlation between two random graphs called *total correlation* is introduced; it is neatly partitioned into an inter-graph contribution (the “edge correlation” that had been the previous focus in the literature) and a novel intra-graph contribution. Furthermore, they introduce a statistic called *alignment strength*, which is 1 minus a normalized count of the number of disagreements in an optimal/true graph match; they prove under mild conditions that alignment strength is a strongly consistent estimator of total correlation. Experimental results in Fishkind et al. (2019b) suggest that the matchability phase transition, as well as the complexity of the problem, is a function of this more nuanced total correlation rather than simply the cross-graph edge correlation/edge sparsity that had been the previous focus in the literature.

Analyses mining the matchability phase transition in the literature that also have considered similarity across generative network models beyond simple sparsity have thus far focused on simple community-structured network models (Onaran et al. 2016; Shirani et al. 2018; Lyzinski 2018), or have proceeded by removing the heterogeneous within-graph model information and simply using the across graph edge correlation (Lyzinski and Sussman 2020). Recently, there have been numerous papers in the literature at the interface between algorithm development and mining matchability phase-transitions; see, for instance, Barak et al. (2019), Mossel and Xu (2020), Ding et al. (2020). A common theme of many of these results is that, under assumptions on the across graph edge-correlation and network sparsity, algorithms are designed to efficiently (or approximately efficiently) match graphs with corresponding theoretical guarantees on the performance of the algorithms in recovering the latent alignment.

However, the question remains how a practitioner knows in practice whether or not a graph matching has successfully recovered the truth. This issue is not resolved by asymptotic analysis with hidden constants. Nor, in general, are the underlying parameters known to the practitioner. It seems that the graph alignment statistic is a very natural metric to use in deciding if the truth is found. Unfortunately, when there is an absence of signal, an optimal (or approximately optimal) graph matching will find spurious and random alignment strength due to chance. Indeed, this meaningless alignment strength can be high and misleading. How do we gauge whether or not it is high enough to signal that truth is found?

After formally defining seeded graph matching and alignment strength in Sect. “[Overview: seeded graph matching, alignment strength](#)” and defining the correlated Bernoulli

random graph model (and attendant parameters) in Sect. “[The correlated Bernoulli random graph model](#)”, we then address this issue with our Phantom Alignment Strength Conjecture in Sect. “[Phantom alignment strength conjecture, consequences](#)”, and in the ensuing discussion in Sect. “[Phantom alignment strength conjecture, consequences](#)”. Then, in Sect. “[Empirical evidence in favor of the phantom alignment strength conjecture](#)”, we present empirical evidence for the conjecture using synthetic and real data, and comparing to theoretical results; Sect. “[Empirical evidence in favor of the phantom alignment strength conjecture](#)” begins with a thorough summary. This is followed in Sect. “[Notable mentions and future directions, plus caveats](#)” by notable mentions, and future directions.

Overview: seeded graph matching, alignment strength

In the seeded graph matching setting, we are given two simple graphs, say they are $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$, such that $|V_1| = |V_2|$, denote the number of vertices $n := |V_1|$. Let Π denote the set of all bijections $V_1 \rightarrow V_2$. It is usually understood that there exists a “true” bijection $\varphi^* \in \Pi$ which represents a natural correspondence between the vertices in V_1 and the vertices in V_2 ; for example, V_1 and V_2 might be the same people, with E_1 indicating which pairs exchanged emails and E_2 indicating pairs that communicated in a different medium. Or G_1 may be the electrical connectome (brain graph) of a worm and G_2 might be the chemical connectome of the same worm, both graphs sharing the same vertex set of neurons. The vertex set V_1 is partitioned into two disjoint sets, \mathcal{S} “seeds” (possibly empty) and \mathcal{N} “nonseeds,” denote $s := |\mathcal{S}|$ and $n := |\mathcal{N}|$. (When $s = 0$ this is the conventional graph matching problem.) The graphs G_1 and G_2 are observed, and the values of φ^* are observed on the set of seeds \mathcal{S} , however the values of φ^* are not observed on the nonseeds \mathcal{N} , and one of several important tasks is to estimate φ^* .

Let $\Pi^{\mathcal{S}}$ denote the set of all bijections $V_1 \rightarrow V_2$ that agree with φ^* on the seeds \mathcal{S} . For any $\varphi \in \Pi^{\mathcal{S}}$, its *match ratio* is defined to be $\frac{1}{n} |\{v \in \mathcal{N} : \varphi(v) = \varphi^*(v)\}|$, i.e. the fraction of the nonseeds that are correctly matched by φ . (It is common to multiply the match ratio by 100 to express it as a percentage.)

For any set V , let $\binom{V}{2}$ denote the set of two-element subsets of V ; for each $i = 1, 2$ and any $\{u, v\} \in \binom{V_i}{2}$ let $u \sim_{G_i} v$ and $u \not\sim_{G_i} v$ denote adjacency and, respectively, nonadjacency of u and v in G_i . Next, let $\mathbf{1}$ denote the indicator function for its subscript. Given any $\varphi \in \Pi$, we define the *full number of disagreements through φ* to be

$$\mathcal{D}'(\varphi) := \sum_{\{u,v\} \in \binom{V_1}{2}} \left(\mathbf{1}_{[u \sim_{G_1} v] \wedge [\varphi(u) \not\sim_{G_2} \varphi(v)]} + \mathbf{1}_{[u \not\sim_{G_1} v] \wedge [\varphi(u) \sim_{G_2} \varphi(v)]} \right) \quad (1)$$

and, given any $\varphi \in \Pi^{\mathcal{S}}$, we define the *restricted number of disagreements through φ* to be

$$\mathcal{D}(\varphi) := \sum_{\{u,v\} \in \binom{\mathcal{N}}{2}} \left(\mathbf{1}_{[u \sim_{G_1} v] \wedge [\varphi(u) \not\sim_{G_2} \varphi(v)]} + \mathbf{1}_{[u \not\sim_{G_1} v] \wedge [\varphi(u) \sim_{G_2} \varphi(v)]} \right). \quad (2)$$

The seeded graph matching problem is to find

$$\hat{\varphi} \in \arg \min_{\varphi \in \Pi^{\mathcal{S}}} \mathcal{D}'(\varphi), \quad (3)$$

and the idea is that $\hat{\varphi}$ is an estimate for the true bijection φ^* . Unfortunately, except in the smallest instances, computing $\hat{\varphi}$ is intractable. A state-of-the-art algorithm SGM from (Fishkind et al. 2019a) is commonly used to approximately solve the optimization problem in (3), and we denote its output $\hat{\varphi}_{SGM}$ ($\in \Pi^S$), and it is an approximation of $\hat{\varphi}$ and, hence, an approximation of φ^* . For any $\varphi \in \Pi^S$, the *full alignment strength* $\text{str}'(\varphi)$ and the *restricted alignment strength* $\text{str}(\varphi)$ are defined as

$$\text{str}'(\varphi) := 1 - \frac{\mathcal{D}'(\varphi)}{\frac{1}{n!} \sum_{\phi \in \Pi} \mathcal{D}'(\phi)} \quad \text{and} \quad \text{str}(\varphi) := 1 - \frac{\mathcal{D}(\varphi)}{\frac{1}{n!} \sum_{\phi \in \Pi^S} \mathcal{D}(\phi)}. \quad (4)$$

Although the denominators of (4) have exponentially many summands, alignment strength is easily computed as follows. For $i = 1, 2$, define the full density of G_i as $\mathfrak{d}'G_i := \frac{|E_i|}{\binom{n}{2}}$ and the restricted density of G_i as $\mathfrak{d}G_i$ = the number of edges of G_i induced by \mathcal{N} , divided by $\binom{n}{2}$. It holds that

$$\begin{aligned} \text{str}'(\varphi) &= 1 - \frac{\mathcal{D}'(\varphi) / \binom{n}{2}}{\mathfrak{d}'G_1(1 - \mathfrak{d}'G_2) + (1 - \mathfrak{d}'G_1)\mathfrak{d}'G_2} \quad \text{and} \\ \text{str}(\varphi) &= 1 - \frac{\mathcal{D}(\varphi) / \binom{n}{2}}{\mathfrak{d}G_1(1 - \mathfrak{d}G_2) + (1 - \mathfrak{d}G_1)\mathfrak{d}G_2}; \end{aligned} \quad (5)$$

see Fishkind et al. (2019b) for the derivation of (5) from (4).

The importance of alignment strength to a practitioner is twofold:

First, the alignment strength of φ^* (and its proxies $\hat{\varphi}$ and $\hat{\varphi}_{SGM}$) may be thought of as a measure of how similar the structure of the graphs G_1 and G_2 are through the “true” bijection; indeed, if the number of disagreements under φ^* [and its proxies $\hat{\varphi}$ and $\hat{\varphi}_{SGM}$] is about equal to the average over all bijections then its alignment strength is near 0 (as clearly seen from the definition in (4)) and, at the other extreme, if φ^* (and its proxies $\hat{\varphi}$ and $\hat{\varphi}_{SGM}$) is nearly an isomorphism between G_1 and G_2 then its alignment strength is near 1. It was proven in Fishkind et al. (2019b) that the full alignment strength of the “true” bijection $\text{str}'(\varphi^*)$ is a strongly consistent estimator of ϱ_T , which is a parameter called the *total correlation* between the two graphs G_1 and G_2 , defined in Sect. “The correlated Bernoulli random graph model”.

Another way that alignment strength is of much importance to a practitioner is in providing confidence that $\hat{\varphi}_{SGM}$ or $\hat{\varphi}$ is a good estimate of φ^* , the “truth.” If $\text{str}(\hat{\varphi}_{SGM})$ or $\text{str}(\hat{\varphi})$ is high enough then we may be confident that a meaningful match capturing similar graph structure has been found, and therefore $\hat{\varphi}_{SGM}$ or $\hat{\varphi}$ is approximately or exactly φ^* . But, how high is high enough?

Indeed, these issues in the use of alignment strength become vastly more complicated by the possibility of *phantom alignment strength*. This is a phenomenon that occurs when, in the presence of weak signal, meaningless matchings have many fewer disagreements than average (sometimes very substantially fewer) due to random noise, and $\hat{\varphi}$ and/or $\hat{\varphi}_{SGM}$ is one of these meaningless matchings—optimal in the optimization problem, but meaningless as estimates of φ^* . Indeed, the alignment strength of $\hat{\varphi}$ and/or $\hat{\varphi}_{SGM}$ may be elevated enough to give the misleading appearance of significance when, in reality, they don’t at all resemble φ^* . This will be illustrated in Sect. “Empirical evidence in favor of the phantom alignment strength conjecture”.

The purpose of this paper is to give a principled, practical means of approaching the decision of what level of alignment strength for $\hat{\varphi}$ and/or $\hat{\varphi}_{SGM}$ indicates that they are a good approximation of φ^* , in which case the alignment strength reflects the amount of meaningful similar structure between G_1 and G_2 —beyond the random similarity between completely unrelated graphs.

(A note on terminology: We define both *full* alignment strength and *restricted* alignment strength since each will end up being important at a different time. The Phantom Alignment Strength Conjecture of Sect. “[Phantom alignment strength conjecture, consequences](#)” requires restricted alignment strength specifically; indeed, since full alignment strength includes the seeds, this would dilute the desired effect, falsifying the conjecture conclusion. However, after we have confidence that our graph matching is the true matching, it is then full alignment strength that will be a better estimator of total correlation introduced in Sect. “[The correlated Bernoulli random graph model](#)”).

The correlated Bernoulli random graph model

Definition 1 Given positive integer n , vertex set V such that $|V| = n$, the parameters of the *correlated Bernoulli random graph model* are Bernoulli parameters $p_{\{u,v\}} \in [0, 1]$ for each $\{u, v\} \in \binom{V}{2}$, and an edge correlation parameter $\varrho_e \in [0, 1]$. The pair of random graphs (G_1, G_2) have a *correlated Bernoulli random graph distribution* when as follows: G_1 and G_2 each have vertex set V . For each $\{u, v\} \in \binom{V}{2}$, and each $i = 1, 2$, the probability of $u \sim_{G_i} v$ is the Bernoulli parameter $p_{\{u,v\}}$, and the Pearson correlation for random variables $\mathbf{1}_{v \sim_{G_1} w}$ and $\mathbf{1}_{v \sim_{G_2} w}$ is the edge correlation parameter ϱ_e . Other than these dependencies, the rest of the adjacencies are independent.

The distribution of the pair of random graphs G_1, G_2 is determined by the above (see Fishkind et al. 2019b). Of course, the identity function is the “true” matching φ^* between G_1 and G_2 .

(If the Bernoulli parameters are all equal, then the random graphs G_1 and G_2 are each said to be Erdos–Renyi, so the correlated Erdos–Renyi random graph model is a special case of the correlated Bernoulli random graph model.)

Important functions of the model parameters are as follows. The Bernoulli mean and Bernoulli variance are, respectively, defined as

$$\mu := \frac{\sum_{\{u,v\} \in \binom{V}{2}} p_{\{u,v\}}}{\binom{n}{2}}, \quad \sigma^2 := \frac{\sum_{\{u,v\} \in \binom{V}{2}} (p_{\{u,v\}} - \mu)^2}{\binom{n}{2}}.$$

Assume that μ is not equal to 0 nor 1. The *heterogeneity correlation* is defined in Fishkind et al. (2019b) as

$$\varrho_h := \frac{\sigma^2}{\mu(1 - \mu)}; \tag{6}$$

it is in the unit interval $[0, 1]$; see Fishkind et al. (2019b). Also pointed out in Fishkind et al. (2019b) is that ϱ_h is 0 if and only if all Bernoulli parameters are equal (i.e. the graphs are Erdos–Renyi) and ϱ_h is 1 if and only if all Bernoulli parameters are $\{0, 1\}$ -valued. In

particular, if ϱ_h is 1 then G_1 and G_2 are almost surely isomorphic. The *total correlation* ϱ_T is defined in Fishkind et al. (2019b) to satisfy the relationship

$$(1 - \varrho_T) = (1 - \varrho_h)(1 - \varrho_e). \quad (7)$$

In the following key result, Theorem 1, which was proved in Fishkind et al. (2019b), let us consider a probability space that incorporates correlated Bernoulli random graph distributions for each of the number of vertices $n = 1, 2, 3, \dots$. Thus, the parameters are functions of n , but to prevent notation clutter we omit notating the dependence on n . The symbol $\xrightarrow{a.s.}$ denotes almost sure convergence.

Theorem 1 *Suppose μ is bounded away from 0 and 1, over all n . Then it holds that $\text{str}'(\varphi^*) - \varrho_T \xrightarrow{a.s.} 0$.*

Theorem 1 together with Eq. 7 shows that the alignment strength of the true bijection captures (asymptotically) an underlying correlation between the random graphs that can be neatly (and symmetrically, per Eq. 7) partitioned into a inter-graph contribution (edge correlation) and an intra-graph contribution (heterogeneity correlation).

Next, instead of considering a sequence of correlated Bernoulli random graphs, let us dig down deeper one probabilistic level. Specifically, suppose that for each $\{u, v\} \in \binom{V}{2}$ there exists an interval-[0, 1]-valued distribution $F_{\{u,v\}}$ such that the Bernoulli parameter $p_{\{u,v\}}$ (in the correlated Bernoulli random graph model) is an independent random variable with distribution $F_{\{u,v\}}$. Denote the mean of this distribution $\mu_{F_{\{u,v\}}}$, denote the variance of this distribution $\sigma_{F_{\{u,v\}}}^2$, and (if we have $\mu_{F_{\{u,v\}}}$ not 0 nor 1) define the heterogeneity correlation of the distribution to be

$$\varrho_{F_{\{u,v\}}} := \frac{\sigma_{F_{\{u,v\}}}^2}{\mu_{F_{\{u,v\}}}(1 - \mu_{F_{\{u,v\}}})}. \quad (8)$$

Theorem 2 *Given an edge correlation parameter $\varrho_e \in [0, 1]$ and, for each $\{u, v\} \in \binom{V}{2}$, given a [0, 1]-valued distribution $F_{\{u,v\}}$ such that the Bernoulli parameter $p_{\{u,v\}}$ is independently distributed as $F_{\{u,v\}}$, then the distribution of the associated correlated Bernoulli random graphs (G_1, G_2) is completely specified by ϱ_e and, for all $\{u, v\} \in \binom{V}{2}$, the values of $\mu_{F_{\{u,v\}}}$ and $\varrho_{F_{\{u,v\}}}$.*

Proof

Consider any $\{u, v\} \in \binom{V}{2}$; the Bernoulli coefficient $p_{\{u,v\}}$, call it X , has distribution $F_{\{u,v\}}$. For any $p \in [0, 1]$, conditioning on $X = p$, the joint probabilities of combinations of u, v adjacency in G_1, G_2 are computed in a straightforward way (see Fishkind et al. 2019b Appendix A) in the table:

	$u \sim_{G_2} v$	$u \not\sim_{G_2} v$	
$u \sim_{G_1} v$	$p^2 + \varrho_e p(1 - p)$	$(1 - \varrho_e)p(1 - p)$	
$u \not\sim_{G_1} v$	$(1 - \varrho_e)p(1 - p)$	$(1 - p)^2 + \varrho_e p(1 - p)$	(9)

Probabilities of these adjacency combinations, relative to the underlying distribution $F_{\{u,v\}}$, are computed by integrating/summing the conditional probabilities (in table) times the density/mass of $F_{\{u,v\}}$, obtaining

$$\begin{aligned}\mathbb{P}[u \sim_{G_1} v \text{ and } u \not\sim_{G_2} v] &= \mathbb{P}[u \not\sim_{G_1} v \text{ and } u \sim_{G_2} v] \\ &= (1 - \varrho_e)(\mathbb{E}X - \mathbb{E}X^2) \\ &= (1 - \varrho_e)(\mathbb{E}X - (\mathbb{E}X)^2 - \mathbb{E}X^2 + (\mathbb{E}X)^2) \\ &= (1 - \varrho_e)[\mu_{F_{\{u,v\}}}(1 - \mu_{F_{\{u,v\}}}) - \sigma_{F_{\{u,v\}}}^2] \\ &= \mu_{F_{\{u,v\}}}(1 - \mu_{F_{\{u,v\}}})(1 - \varrho_e)(1 - \varrho_{F_{\{u,v\}}}).\end{aligned}$$

Then, for each $i = 1, 2$, because $\mathbb{P}[u \sim_{G_i} v] = \mathbb{E}X = \mu_{F_{\{u,v\}}}$ we have all four adjacency combinations as functions of $\mu_{F_{\{u,v\}}}$ and $\varrho_{F_{\{u,v\}}}$. The result follows from the independence across all pairs of vertices. \square

In the Phantom Alignment Strength Conjecture we assume all distributions $F_{\{u,v\}}$ are the same, call the common distribution F . Note that Bernoulli mean μ and heterogeneity correlation ϱ_h are now random variables, and if n is large, then μ and ϱ_h will respectively be good estimators of μ_F and ϱ_F . A very important consequence of Theorem 2 is that the only information that matters regarding F is contained (well-estimated) in the quantities μ and ϱ_h .

Phantom alignment strength conjecture, consequences

In this section, we propose the Phantom Alignment Strength Conjecture, which is the central purpose of this paper. We then discuss its consequences; the conjecture gives us a principled and practical way to decide if we should be convinced that the output of a graph matching algorithm well-approximates the true matching.

Henceforth we use the term *alignment strength* to refer to the restricted alignment strength.

Consider correlated Bernoulli random graphs G_1, G_2 such that there are a “moderate” number n of nonseed vertices (say $n \geq 300$), s seeds (selected discrete uniformly from the $n := n + s$ vertices), and Bernoulli parameters are independently realized from any fixed $[0, 1]$ -valued distribution with moderate mean μ' (say $.05 < \mu' < .95$). The *Phantom Alignment Strength Conjecture* states that, subject to caveats, as discussed in Sect. “[Notable mentions and future directions, plus caveats](#)”, there exists a *phantom alignment strength value* $\hat{q} \equiv \hat{q}(n, s, \mu') \in [0, 1]$ such that $\text{str}(\hat{\varphi})$ has “negligible” variance and is approximately a function of the total correlation ϱ_T and, specifically, it holds that, with “high probability,”

$$\text{str}(\hat{\varphi}) \approx \begin{cases} \varrho_T & \text{if } \varrho_T > \hat{q}; \text{ in which case } \hat{\varphi} = \varphi^* \\ \hat{q} & \text{if } \varrho_T \leq \hat{q}; \text{ in which case } \hat{\varphi} \text{ is “very different from” } \varphi^* \end{cases} \quad (10)$$

Moreover, the conjecture states that, when using the seeded graph matching algorithm SGM of Fishkind et al. (2019a), (given n, s, μ' , as above) then there exists $\hat{q}_{SGM} \equiv \hat{q}_{SGM}(n, s, \mu') \in [0, 1]$ such that $\hat{q}_{SGM} \geq \hat{q}$, and $\text{str}(\hat{\varphi}_{SGM})$ has “negligible” variance and is approximately a function of the total correlation ϱ_T and, specifically, it holds that, with “high probability,”

$$\text{str}(\hat{\varphi}_{SGM}) \approx \begin{cases} \varrho_T & \text{if } \varrho_T > \hat{q}_{SGM}; \text{ in which case } \hat{\varphi}_{SGM} = \varphi^* \\ \hat{q} & \text{if } \varrho_T \leq \hat{q}_{SGM}; \text{ in which case } \hat{\varphi}_{SGM} \text{ is "very different" from } \varphi^* \end{cases} \quad (11)$$

Note that both $\text{str}(\hat{\varphi})$ and $\text{str}(\hat{\varphi}_{SGM})$ are conjectured to be an approximately piecewise linear function of ϱ_T ; two pieces, one piece with slope 0 and one piece with slope 1. However, $\text{str}(\hat{\varphi})$ is continuous and shaped like a hockey stick (see Fig. 2f), whereas for $\text{str}(\hat{\varphi}_{SGM})$ there can be a discontinuity (see Fig. 2b); but the function value of the linear portion with slope 0 is the same for $\text{str}(\hat{\varphi}_{SGM})$ as it is for $\text{str}(\hat{\varphi})$, namely it is the phantom alignment strength value \hat{q} .

There are important consequences of the Phantom Alignment Strength Conjecture for the practitioner. Suppose that a practitioner has two particular graphs G_1, G_2 with n nonseed vertices and s seeds that can be considered as realized from a correlated Bernoulli random graph model, and the practitioner wants to seeded graph match them, computing $\hat{\varphi}_{SGM}$ as an approximation of the true matching φ^* . How can the practitioner tell if $\hat{\varphi}_{SGM}$ is φ^* ? This conjecture provides a principled, practical mechanism. The practitioner should realize two independent Erdos–Renyi graphs H_1 and H_2 with n nonseed vertices, s seeds, and adjacency probability parameter p equal to the combined density of G_1 and G_2 . Then use SGM to seeded graph match H_1 and H_2 , and the alignment strength of the bijection (between H_1 and H_2) is approximately $\hat{q} \equiv \hat{q}(n, s, \mu)$, since the total correlation in generating H_1 and H_2 is 0, by design. Then, when subsequently seeded graph matching G_1 and G_2 , if $\text{str}(\hat{\varphi}_{SGM})$ is greater than some predetermined and fixed $\epsilon > 0$ above \hat{q} , then that would indicate that $\hat{\varphi}_{SGM} = \varphi^*$ and, if $\text{str}(\hat{\varphi}_{SGM})$ is less than this, then there is no confidence that $\hat{\varphi}_{SGM}$ is φ^* . Moreover, in the former case the practitioner can have confidence in approximating $\text{str}(\hat{\varphi}_{SGM}) \approx \varrho_T$, and in the latter case there wouldn't be confidence in this approximation. (In the former case, note that the full alignment strength $\text{str}'(\hat{\varphi}_{SGM})$ would then be an even better estimate of ϱ_T .)

(If some of the model assumptions are violated and the Bernoulli mean of G_1 may be different from G_2 , then it may be better not to combine their densities, but rather to realize H_1 and H_2 as Erdos–Renyi graphs with respective adjacency parameter equal to their respective densities.)

Empirical evidence in favor of the phantom alignment strength conjecture

In this section we provide empirical evidence for the Phantom Alignment Strength Conjecture.

A summary is as follows:

We begin in Sect. “Of hockey sticks and phantom alignment strength” with a scale small enough (n is just on the order of tens) to solve seeded graph matching and attain optimality. Although the Phantom Alignment Strength Conjecture does not apply because n is so small, we nonetheless see many ingredients of the conjecture. Then, in Sect. “Of hockey sticks and broken hockey sticks”, we use synthetic data on a scale for the conjecture to be applicable, and we empirically demonstrate the conjecture for many types of Bernoulli parameter distributions; unimodal, bimodal, symmetric, skewed, etc. The SGM algorithm is employed for seeded graph matching, since exact optimality is unattainable in practice.

In Sect. “[Phantom alignment strength versus theoretical matchability threshold](#)”, the alignment strength of completely uncorrelated Erdos–Renyi graphs (graph matched with SGM, using no seeds), taken as a function of n , is empirically demonstrated to be the same order of growth (in terms of n) as the theoretical bound for matchability (as a function of n), which suggests that the two quantities are the same, in excellent accordance with the conjecture.

Then, in Sect. “[Block settings](#)”, we observe that when there is block structure and differing distributions for the Bernoulli parameters by block (thus the conjecture hypotheses are not adhered to) then the conjecture’s claims may fail to hold, to some degree. Nonetheless, there is still a phantom alignment strength that allows for a procedure similar to what we recommend in Sect. “[Phantom alignment strength conjecture, consequences](#)” to be successfully used for deciding when alignment strength is significant enough to indicate that the seeded graph matching has found the truth.

Real data is then used for demonstration in Sects. “[Real data; matching graphs to noisy renditions](#)” and “[Real data; matching same objects under different modalities](#)”.

Specifically, in Sect. “[Real data; matching graphs to noisy renditions](#)”, we use a human connectome at many different resolution levels, and graph match it to a manually noised copy of itself.

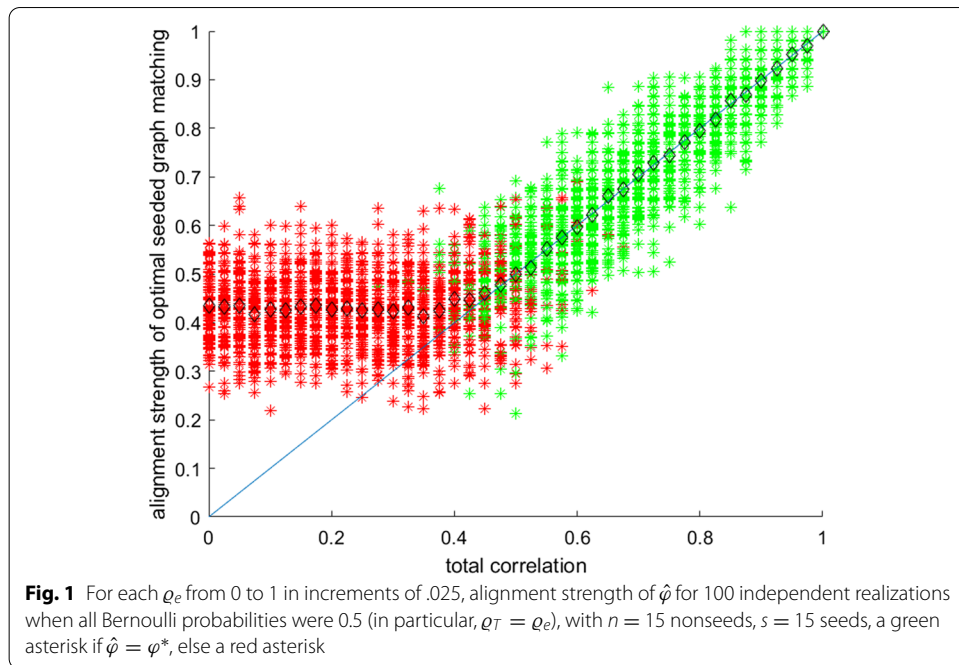
Then, in Sect. “[Real data; matching same objects under different modalities](#)”, we consider several pairs of real-data graphs (titled Wikipdeia, Enron, and C Elegans) whose vertices are the same objects, and the adjacencies in each pair of graphs represent relationships between the objects across two different modalities.

All of these experiments serve as strong empirical evidence for the Phantom Alignment Strength Conjecture, and motivate its use.

Of hockey sticks and phantom alignment strength

We begin with an experiment in which the value of n is well below what is required in the statement of the Phantom Alignment Strength Conjecture. However, n is small enough here to enable us to compute $\hat{\varphi}$ exactly, using the integer programming formulation from Fishkind et al. (2019b). We will be able to see many features of the Phantom Alignment Strength Conjecture, and we will also see that phantom alignment strength is not just an artifact of the SGM algorithm.

For each value of ϱ_e from 0 to 1 in increments of .025, we did 100 independent repetitions of the following experiment. We realized a pair of correlated Bernoulli random graphs on $n = 30$ vertices with edge correlation ϱ_e and, for each pair of vertices, the associated Bernoulli parameter was 0.5. (In particular, the graphs are correlated Erdos–Renyi.) Since here $\sigma^2 = 0$, we have that $\varrho_h = 0$, and thus $\varrho_T = \varrho_e$. We discrete uniform randomly chose $s = 15$ seeds, so there were $n = 15$ nonseeds. For each experiment, we solved the seeded graph matching problem to optimality (indeed, $n = 15$ is small enough to do so), obtaining $\hat{\varphi}$. If it happened that $\hat{\varphi} = \varphi^*$ then we plotted a green asterisk in Fig. 1 for the resulting alignment strength $\text{str}(\hat{\varphi})$ against the total correlation ϱ_T and, if $\hat{\varphi} \neq \varphi^*$, we plotted a red asterisk for the resulting alignment strength $\text{str}(\hat{\varphi})$ against the total correlation ϱ_T . The black diamonds in Fig. 1 are the mean alignment strengths for the 100 repetitions, plotted for each value of ϱ_e .



It is readily seen from Fig. 1 that the variance for the alignment strength of $\hat{\varphi}$ is quite high, which is reason to not formulate the Phantom Alignment Strength Conjecture until n is much larger. Other than this, observe that if we substitute “mean of the alignment strength of $\hat{\varphi}$ ” into the conjecture in place of “alignment strength of $\hat{\varphi}$ ” then the conjecture would hold here. Indeed, when $\varrho_T > \approx 0.44 \equiv \hat{q}$ we very generally had that $\hat{\varphi} = \varphi^*$, and when $\varrho_T \leq \approx 0.44$ we very generally had that $\hat{\varphi} \neq \varphi^*$. (This boundary is not sharp, but is close.) Also, note that when $\varrho_T > \approx 0.44$, the mean of the alignment strength was approximately equal to ϱ_T . Furthermore, when $\varrho_T \leq \approx 0.44$, we see that the (mean) alignment strength of $\hat{\varphi}$ is the phantom alignment strength (mean) of ≈ 0.44 . Indeed, in this latter case, the alignment strength of $\hat{\varphi}$ is a misleading high value, and is not meaningful.

Of hockey sticks and broken hockey sticks

In this section, we use synthetic data that meets the hypotheses of the Phantom Alignment Strength Conjecture. Our setup was as follows. We chose the number of nonseeds to be $n = 1000$, and we repeated an experiment for **all combinations** of the following:

- Each pair of Beta distribution parameters α, β listed in the following table:

	α	β
Pair A	1	1
Pair B	0.5	0.5
Pair C	2	2
Pair D	5	1
Pair E	2	5

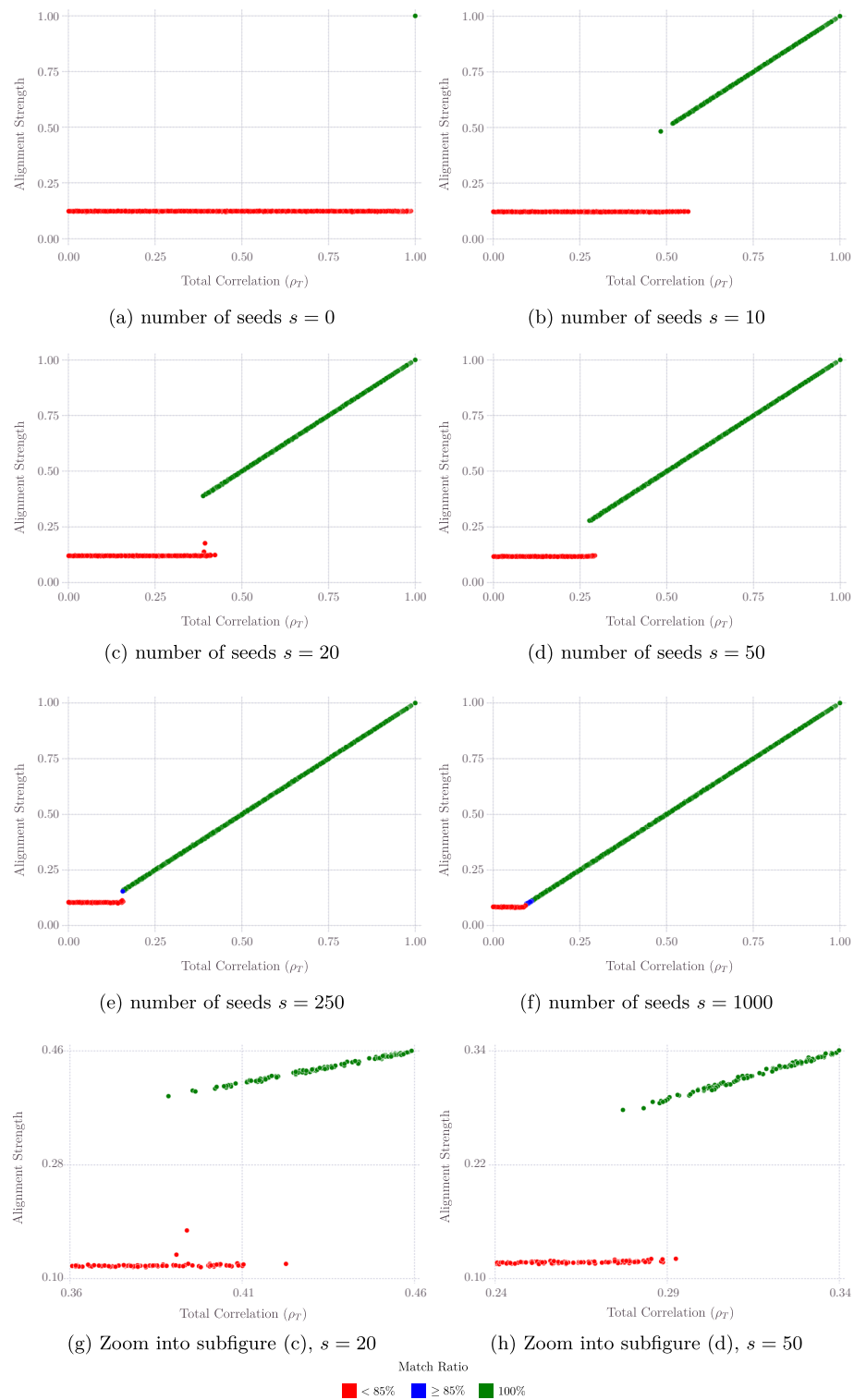


Fig. 2 Alignment strength $\text{str}(\hat{\varphi}_{SGM})$ plotted against total correlation ρ_T for the synthetic data experiments in Sect. “Of hockey sticks and broken hockey sticks”, separated according to the number of seeds s . The number of nonseeds was $n = 1000$, and only the case of $\mu' = 0.5$ is shown here. Match ratio of each experiment is color coded green, blue, or red according to the legend above. Subfigures **g** and **h** are zooms into subfigures **c** and **d**, to increase the granularity so that the thresholding is better seen

- Each μ' (mean of the scaled/translated Beta distribution) from .1 to .9 in increments of .1,
- Each number of seeds $s = 0, 10, 20, 50, 250, 1000$,
- Each value of edge correlation ϱ_e from 0 to 1 in increments of 0.025,
- Each value of δ from 0 to $\delta_{max} := \min\{\frac{\alpha+\beta}{\alpha}\mu', \frac{\alpha+\beta}{\beta}(1-\mu')\}$ in increments of $\frac{1}{10}\delta_{max}$.

For each combination of the above, we realized a pair of correlated Bernoulli random graphs on $n + s$ vertices, with edge correlation ϱ_e and, for each pair of vertices, the associated Bernoulli parameter was independently realized from the distribution $\delta \cdot \text{Beta}(\alpha, \beta) + \mu' - \delta \frac{\alpha}{\alpha+\beta}$. Note that

- The distribution $\delta \cdot \text{Beta}(\alpha, \beta) + \mu' - \delta \frac{\alpha}{\alpha+\beta}$ has support interval of length δ , has mean μ' , and the support interval is contained in the interval $[0, 1]$.
- The distribution $\delta \cdot \text{Beta}(\alpha, \beta) + \mu' - \delta \frac{\alpha}{\alpha+\beta}$ is uniform when α, β is 1, 1, and is bimodal when α, β is 0.5, 0.5, is symmetric unimodal when α, β is 2, 2, and is skewed in the other two cases, in different directions, one where the mode is an endpoint of the support and one where the mode is interior of the support.
- The Bernoulli mean μ is approximately μ' , since $\binom{n+s}{2}$ is very large for these purposes.

The s seeds were chosen discrete uniform randomly from the $n + s$ vertices, and we computed $\hat{\varphi}_{SGM}$ via the SGM algorithm for seeded graph matching. In Fig. 2 we plotted alignment strength $\text{str}(\hat{\varphi}_{SGM})$ against total correlation ϱ_T for all of the pairs of graphs generated in the case where $\mu' = 0.5$, in different subfigures for the different values of $s = 0, 10, 20, 50, 250, 1000$; green dots indicate when $\hat{\varphi}_{SGM} = \varphi^*$, blue and red dots indicate when $\hat{\varphi}_{SGM} \neq \varphi^*$, blue when $\hat{\varphi}_{SGM}$ agreed with φ^* on at least 85% of the nonseeded vertices (i.e. “match ratio $\geq 85\%$ ”), and red when $\hat{\varphi}_{SGM}$ agreed with φ^* on less than 85% of the nonseeded vertices.

Note that in Fig. 2, each of (a)–(f) are plots of 2255 points, each point represented with a filled circle, and the crowding of the points makes them resemble lines; so, in Fig. 2, we also included (g) and (h), which are zooms of a portion of (c) and (d), respectively. With the increased granularity in (g) and (h), we see that if we ignore some outlier red and green dots, then there is a better defined transition from red to green than would appear in (c) and (d).

The Phantom Alignment Strength Conjecture is well motivated by the results illustrated in Fig. 2. In particular, alignment strength $\text{str}(\hat{\varphi}_{SGM})$ exhibits very low variance and is approximately a piecewise-linear function of total correlation ϱ_T . There appears to be a critical value \hat{q}_{SGM} , dependent on the number of seeds s in these experiments, for which the following holds. When total correlation ϱ_T is above \hat{q}_{SGM} then $\hat{\varphi}_{SGM} = \varphi^*$ and $\text{str}(\hat{\varphi}_{SGM}) \approx \varrho_T$, and when total correlation ϱ_T is below \hat{q}_{SGM} then $\hat{\varphi}_{SGM} \neq \varphi^*$, evidenced by $\text{str}(\hat{\varphi}_{SGM}) \not\approx \varrho_T$, and $\text{str}(\hat{\varphi}_{SGM})$ is constant—at a phantom alignment strength level. When there are enough seeds, we see that the two pieces of the function join to become continuous, suggesting that $\hat{\varphi}_{SGM} = \hat{\varphi}$ is then achieved for all ϱ_T , and the value of \hat{q}_{SGM} is then \hat{q} .

Also note that the five different Beta distributions from which Bernoulli parameters were realized (the five pairs of Beta parameters labelled A, B, C, D, E) in these

experiments were collected into each of the figures of Fig. 2, and the experiment results for these different distributions are indistinguishable from each other in the figures, in accordance with Theorem 2, and reflected in the Phantom Alignment Strength Conjecture claim that the phantom alignment strength is just a function of n, s, μ' , and that it isn't relevant what distribution is used to obtain the Bernoulli parameters.

Also note the phase transition from matchable to non-matchable which takes place when Q_T gets to \hat{q}_{SGM} , and this phase transition becomes better and better defined as the number of seeds goes up.

For the other values of μ' , the figures exhibited the same overall type of structure, although the phantom alignment strength values were different. In the interest of space, we only present here the $\mu' = 0.5$ experiment figures.

Phantom alignment strength versus theoretical matchability threshold

Among other assertions, the Phantom Alignment Strength Conjecture asserts, under conditions, that the alignment strength $\text{str}(\hat{\varphi}_{SGM})$ when $Q_T = 0$, called the “phantom alignment strength,” is equal to the total correlation threshold for matchability of exact seeded graph matching (i.e. the particular value such that $\hat{\varphi} = \varphi^*$ or not according as Q_T is greater than this value or not); indeed, we have denoted this common quantity \hat{q} . In this section, we will compare alignment strength $\text{str}(\hat{\varphi}_{SGM})$ when $Q_T = 0$ to the matchability threshold proved in Lyzinski et al. (2014).

Consider a probability space with a sequence of correlated Bernoulli random graphs for each of the number of vertices $n \equiv n = 1, 2, 3, \dots$, with $s = 0$ seeds and all Bernoulli parameters equal to a fixed value p (ie correlated Erdos–Renyi random graphs). When we say that a sequence of events happens “almost always” we mean that, with probability 1, all but a finite number of the events occur. The following result was stated and proved in Lyzinski et al. (2014); although stated there in terms of Q_e , we write Q_T instead, since here, where $Q_h = 0$, we have that $Q_T = Q_e$.

Theorem 3 *There exists positive, real valued, fixed constants c_1, c_2 such that if $Q_T \geq c_1 \sqrt{\frac{\log n}{n}}$ then almost always $\hat{\varphi} = \varphi^*$, and if $Q_T \leq c_2 \sqrt{\frac{\log n}{n}}$ then $\lim_{n \rightarrow \infty} \mathbb{E}|\{\varphi \in \Pi : D'(\varphi) < D'(\varphi^*)\}| = \infty$.*

For each value of $p = .05, .1, .2, .3, .4, .5$, and each of 500 values of n between 500 and 4000, (as mentioned, $s = 0$) we plotted realizations of alignment strength $\text{str}(\hat{\varphi}_{SGM})$ vs the value of n , for uncorrelated ($Q_e = 0$) pairs of random Bernoulli (Erdos–Renyi) graphs where each Bernoulli parameter is p , hence $Q_T = 0$ (since $Q_e = 0, Q_h = 0$). Figure 3 shows the plots for $p = 0.05, 0.1, 0.5$.

Then, for each p , we fit the associated points to a curve $f_p(n) := d_p + c_p \sqrt{\frac{\log n}{n}}$ for real numbers c_p and d_p ; the values of d_p and c_p are given in Table 1, and f_p is also drawn in Fig. 3. For each value of p , note the near-perfect fit of f_p to the associated points plotted in Fig. 3, and note that the value of d_p is close to zero.

Indeed, this suggests, as conjectured in the Phantom Alignment Strength Conjecture, that the phantom alignment strength (ie $\text{str}(\hat{\varphi}_{SGM})$ when $Q_T = 0$) exists as a value \hat{q} which coincides with the amount of total correlation needed for $\hat{\varphi} = \varphi^*$.

Table 1 Values of the constants in $f_p(n) := d_p + c_p \sqrt{\frac{\log n}{n}}$

ρ	d_p	c_p
0.05	-0.021	2.19
0.1	-0.010	1.80
0.2	-0.003	1.58
0.3	-0.001	1.51
0.4	0.000	1.48
0.5	0.000	1.47

Block settings

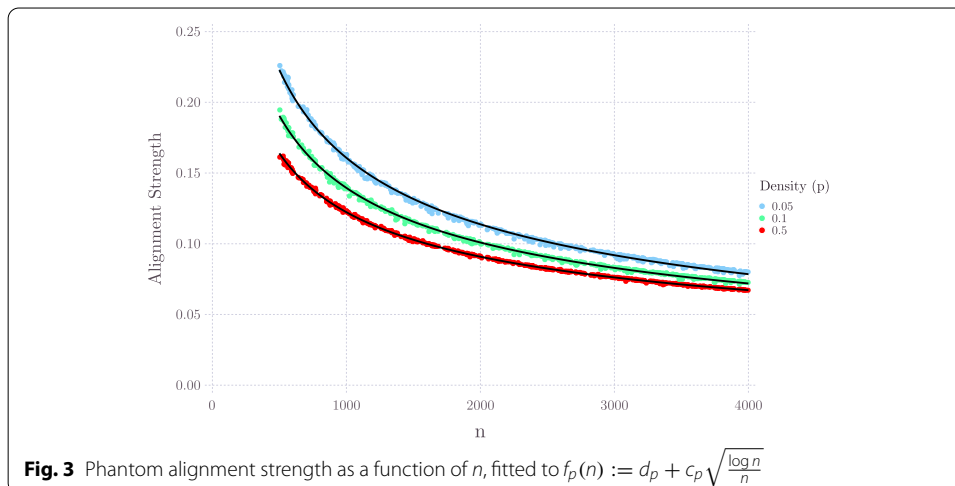
The setting of the Phantom Alignment Strength Conjecture in Sect. “[Phantom alignment strength conjecture, consequences](#)” was specifically concerning correlated Bernoulli random graphs G_1, G_2 such that there are n nonseed vertices, s seed vertices (selected discrete uniformly from the $n := n + s$ vertices), and Bernoulli parameters for each pair of vertices are selected independently from any fixed distribution with mean μ' .

Let us consider a block setting, which differs from the above in that there is a positive integer K , and the vertex set V is first randomly partitioned into K blocks B_1, B_2, \dots, B_K as follows: There is a given probability vector $\pi \in [0, 1]^K$ such that $\sum_{i=1}^K \pi_i = 1$ and each vertex in V is independently placed in block B_i with probability π_i for $i = 1, 2, \dots, K$. Next, suppose there is a unit-interval-valued (ie $[0, 1]$ -valued) distribution $F_{i,j}$ for each $i = 1, 2, \dots, K$ and $j = i, i + 1, \dots, K$ such that, for each $1 \leq i \leq j \leq K$ and each $u \in B_i$ and $v \in B_j$, the Bernoulli parameter $p_{\{u,v\}}$ is independently realized from distribution $F_{i,j}$. Let M be the $K \times K$ symmetric matrix with i, j th entry equal to the mean of distribution $F_{i,j}$.

Similarly to the Phantom Alignment Strength Conjecture, does there exists a phantom alignment strength value $\hat{q} \equiv \hat{q}(n, s, \pi, M) \in [0, 1]$ and also $\hat{q}_{SGM} \equiv \hat{q}_{SGM}(n, s, \pi, M) \in [0, 1]$ whereby Eqs. (10) and (11) hold? This is not so simple.

We consider the following choices for n, s, π , and M :

$$n = 1000 \quad s = 40 \quad \pi = \begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix} \quad M = \begin{bmatrix} 0.3 & 0.4 \\ 0.4 & 0.5 \end{bmatrix}$$



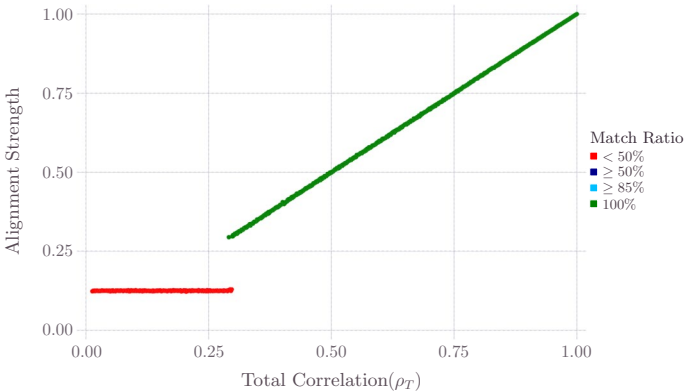


Fig. 4 Experiment A in Sect. “Block settings”; here $F_{1,1}, F_{1,2}, F_{2,2}$ are resp. point mass at 0.3, 0.4, 0.5

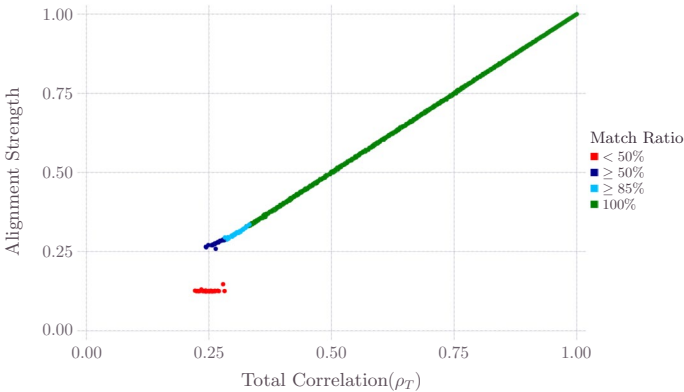


Fig. 5 Experiment B in Sect. “Block settings”; same as Experiment A except that $F_{2,2}$ is uniform [0, 1]

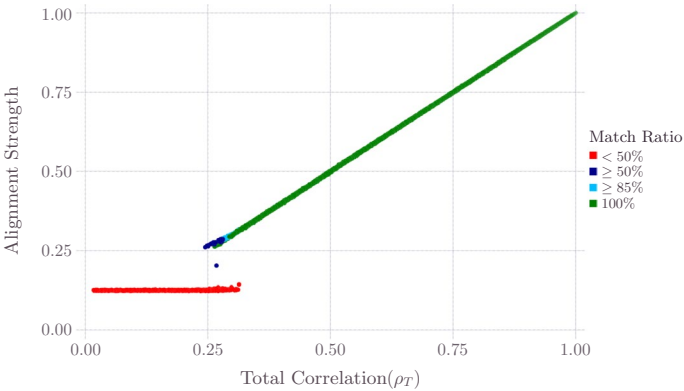


Fig. 6 Experiment C in Sect. “Block settings”; eight different combinations for $F_{1,1}, F_{1,2}, F_{2,2}$

In experiment “A,” we took $F_{1,1}$ to be point mass distribution at 0.3, $F_{1,2}$ to be point mass distribution at 0.4, and $F_{2,2}$ to be point mass distribution at 0.5. For each value of edge correlation ϱ_e from 0 to 1 in increments of 0.001, we realized Bernoulli parameters and then we realized associated correlated Bernoulli random graphs. In Fig. 4, we plotted alignment strength $\text{str}(\hat{\varphi}_{SGM})$ against total correlation ϱ_T ; green dots indicate when $\hat{\varphi}_{SGM} = \varphi^*$, (else) light blue when $\hat{\varphi}_{SGM}$ agreed with φ^* on at least 85% of the nonseeded vertices, (else) dark blue when $\hat{\varphi}_{SGM}$ agreed with φ^* on at least 50%, (else) red when $\hat{\varphi}_{SGM}$ agreed with φ^* on less than 50% of the nonseeded vertices. We then repeated the experiment with the only difference being that $F_{2,2}$ was the uniform distribution on the interval $[0, 1]$, so (n, s, π, M) are same as above; the resulting plot is Fig. 5 (alignment strength $\text{str}(\hat{\varphi}_{SGM})$ vs ϱ_T , same dot color scheme as above). Let us call this Experiment “B.”

Next, we repeated the above experiment for all eight possible combinations of:

$F_{1,1}$ is the uniform distribution on (a) interval $[0.25, 0.35]$ or b) interval $[0, 0.6]$

$F_{1,2}$ is the uniform distribution on (a) interval $[0.35, 0.45]$ or b) interval $[0, 0.8]$

$F_{2,2}$ is the uniform distribution on (a) interval $[0.45, 0.55]$ or b) interval $[0, 1]$

and we superimposed all of the alignment strength vs total correlation plots in Fig. 6 (same dot color scheme as above); we will call this Experiment “C.” Again, the underlying (n, s, π, M) are the same as the previous experiments.

Note that Figs. 4, 5, and 6 (for respective experiments A, B, and C) are not similar, even though they originate from the same values of n, s, π , and M . Thus, the Phantom Alignment Strength Conjecture is not simply extended to the case of nontrivial block structure.

However, also note that when SGM was broadly failing to get the truth in experiments A, B, and C (i.e. the red dots in Figs. 4, 5, and 6), the alignment strength was almost constant, at a value of around 0.12. This suggests a decision procedure (analogous the procedure described in Sect. “Phantom alignment strength conjecture, consequences”) for deciding if G_1, G_2 from an (n, s, π, M) -block model are graph matched with some truth. The procedure would be to realize H_1 and H_2 as correlated Bernoulli random graphs where $\varrho_e = 0$, where the $n + s$ vertices are apportioned to the blocks in proportion to π , and where, for every pair of vertices, the Bernoulli parameter is taken as the entry of M associated with the block memberships of the two vertices, and then the s seeds are chosen uniformly at random. The alignment strength of the seeded graph match of H_1 to H_2 can then be used as a phantom alignment strength value in the sense that, if the alignment strength of the seeded graph match of G_1 to G_2 is more than some $\epsilon > 0$ greater than this phantom alignment strength value, then we decide that there is at least some truth present in the seeded graph match of G_1 to G_2 .

What made the block structure more complicated? We will next provide some insight. Indeed, Experiment B was constructed in an extreme way in order to cause particular mischief. The value of ϱ_h in Experiment A was approximately .0129, and the value of ϱ_h in Experiment B was approximately .2277; in particular, that is why the value of ϱ_T was never below approximately .22 in Experiment B, as is clear from Fig. 5. However, in Experiment B when $\varrho_e = 0$, all of the vertices in the first block are stochastic twins; they share the same probabilities of adjacency as each other to all of the vertices in the graph, and all adjacencies are collectively independent. Thus the “true” bijection (the identity) has no signal in that case. (One might even say that the

“truth” isn’t very “truthy.”) As such, the total correlation in that case, approximately .2277, does not contribute to matchability vis-a-vis the first block. As positive edge correlation ρ_e is increasingly added in to Experiment B, the first block achieves matchability on the strength of only the edge correlation, and the second block achieves matchability on the strength of edge correlation together with heterogeneity correlation. In this manner, total correlation does not tell a uniform story across all vertices. This is in contrast to the hypotheses of the Phantom Alignment Strength Conjecture (and the setup in the empirical matchability experiments in the paper Fishkind et al. 2019b) where the Bernoulli parameters were realized from one distribution. Note that with Experiment C, there is more variety in ρ_h (for the eight experiments the values of ρ_h ranged from approximately .0161 to approximately .30); there is still some lack of demarcation between matchable and nonmatchable in terms of total correlation, but the situation is improved somewhat from the left tail of the figure, and total correlation has more influence as a unified quantity.

We did additional experiments with other values of (n, s, π, M) and found comparable results to what appears above.

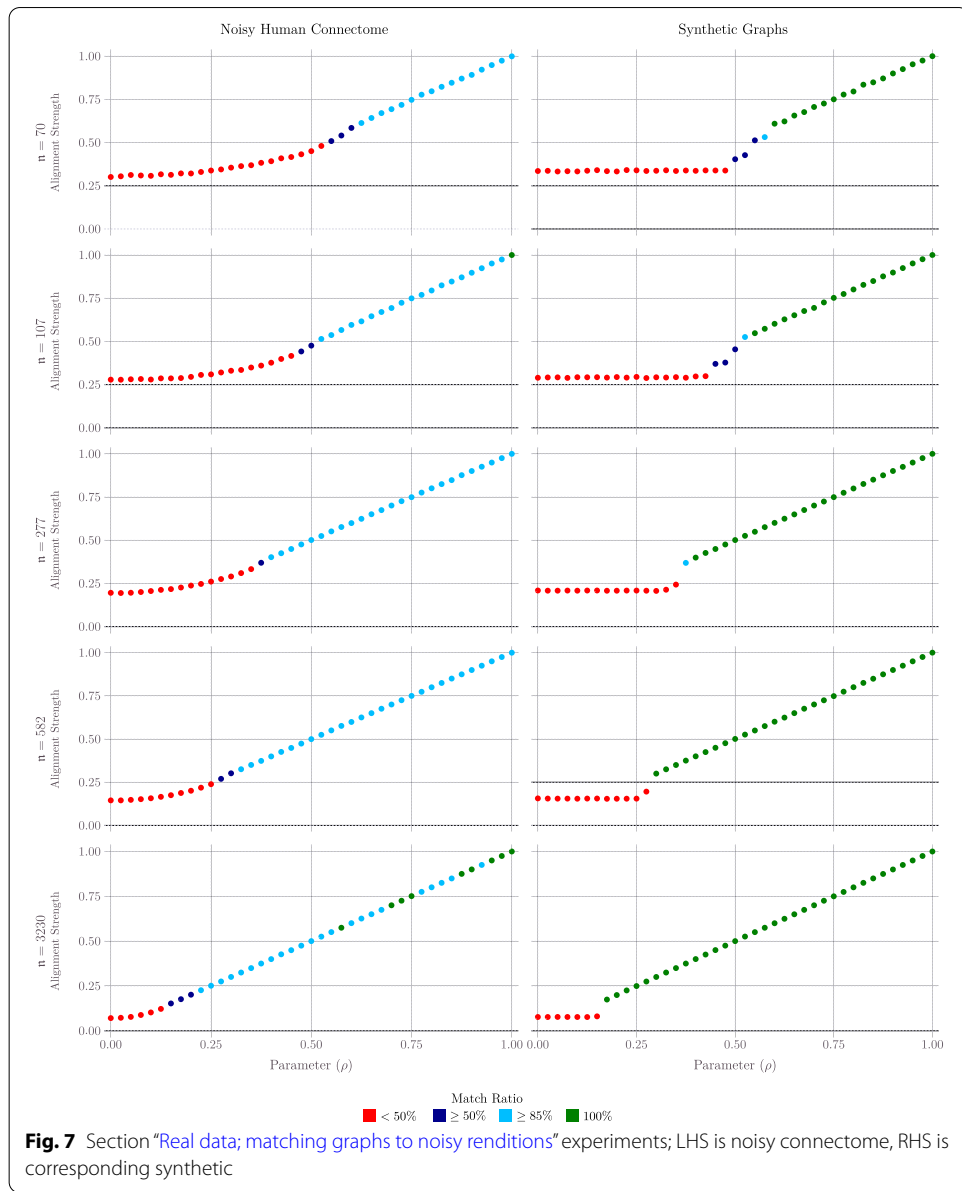
Real data; matching graphs to noisy renditions

Recall that the Phantom Alignment Strength Conjecture is formulated under the assumption that each pair of vertices has a Bernoulli parameter that is a realization of a distribution which is common to all of the pairs of vertices. How realistic is this assumption in practice? And, more to the point of the practitioner, do the conclusions of the conjecture apply to real data, in general?

In this section we consider a human connectome at different resolution levels. (This connectome has been featured in Priebe et al. 2019; Chung et al. 2020.) Diffusion-weighted Magnetic Resonance Imaging (dMRI) brain scans were collected from one hundred and fourteen humans at the Beijing Normal University (Zuo et al. 2014). Fiber tracts, which trace axonal pathways through a three-spatial-dimensional cuboid array of $1 \times 1 \times 1 \text{ mm}^3$ voxels of the dMRI scan, are estimated using the ndmg pipeline (Kiar et al. 2018).

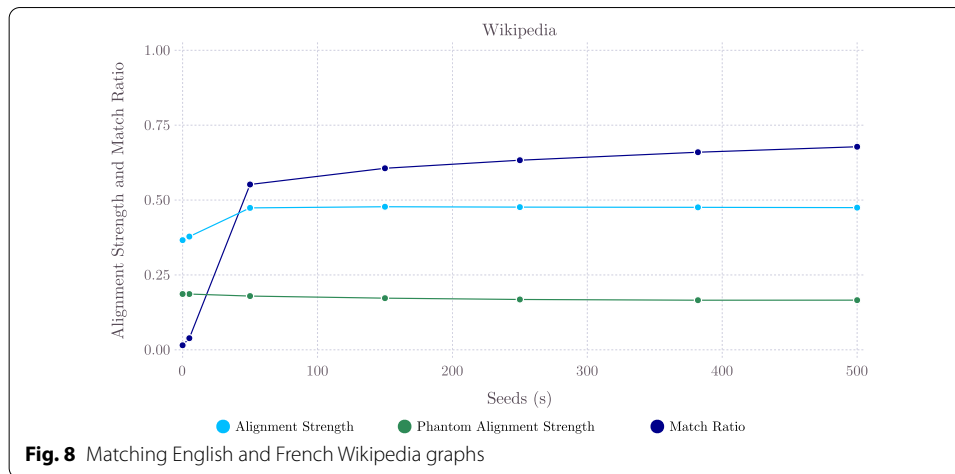
For each value of $n = 70, 107, 277, 582, 3230$, the graph G_n was formed in the following manner. Starting from the original cuboid array of voxels, n equally spaced “contractile” voxels were selected, and each voxel in the array was merged with its nearest contractile voxel (Mhembere et al. 2013); the n such groupings of voxels (centered at their contractile voxel) are the n vertices of the graph G_n . For any two vertices in G_n , we declare them adjacent precisely when there exists a fiber that runs through any voxel of one vertex and also any voxel of the other vertex for any of the one hundred and fourteen individuals.

Given any graph $G = (V, E)$, and also given any noise parameter $\rho \in [0, 1]$, we can instantiate a graph \bar{G} called a ρ -noised rendition of G on the same vertex set V as follows. Denote the density of G by $\bar{\rho}'G := \frac{|E|}{\binom{|V|}{2}}$. First, instantiate an independent Erdos–Renyi graph H on V with Bernoulli parameter $\bar{\rho}'G$; i.e. each pair of vertices is an edge independently of the others with probability $\bar{\rho}'G$. Next, for each pair of vertices $\{u, v\}$, perform an independent Bernoulli trial; with probability ρ set u adjacent/ not adjacent (resp.) to v in



\bar{G} according as u adjacent/ not adjacent (resp.) to v in G , and with probability $1 - \rho$ set u adjacent/ not adjacent (resp.) to v in \bar{G} according as u adjacent/ not adjacent (resp.) to v in H . In this manner, \bar{G} is a mixture of G and noise graph H . When graph matching G to a ρ -noised rendition of G , clearly φ^* is the identity function V to V .

For each of $n = 70, 107, 277, 582, 3230$, we did the following experiment. For each value of the noise parameter ρ from 0 to 1 in increments of .025, we did 20 repetitions of instantiating a ρ -noised rendition of G_n , then seeded graph matched G_n to it using the SGM algorithm after selecting 10% of the n vertices (discrete uniform randomly) as seeds. The mean alignment strength $\text{str}(\hat{\varphi}_{SGM})$ (the mean being over the 20 repetitions) vs noise parameter ρ was plotted in five respective figures (for the five different values of n) in the left side of Fig. 7; green dots indicate when $\hat{\varphi}_{SGM} = \varphi^*$, (else) light blue when



$\hat{\varphi}_{SGM}$ agreed with φ^* on at least 85% of the nonseeded vertices, (else) dark blue when $\hat{\varphi}_{SGM}$ agreed with φ^* on at least 50%, (else) red when $\hat{\varphi}_{SGM}$ agreed with φ^* on less than 50% of the nonseeded vertices.

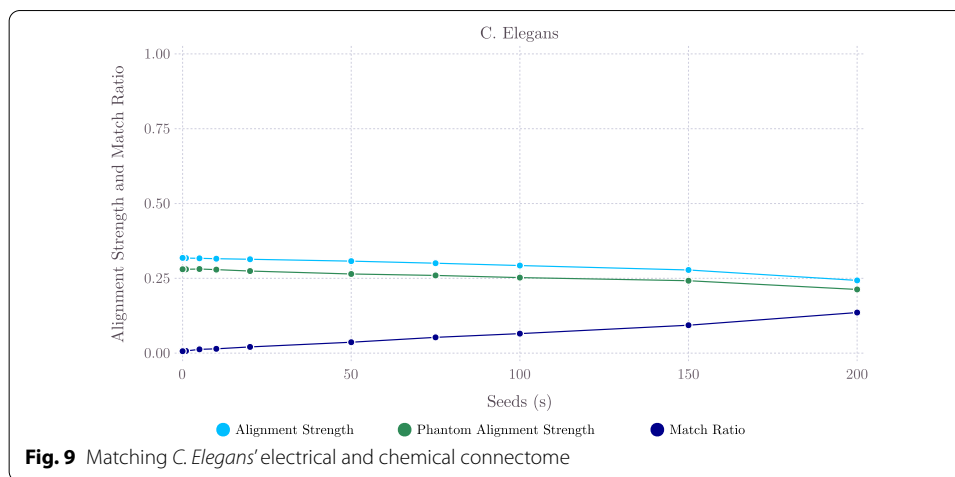
We then repeated the above experiments, with the only difference being that in place of G_n we used an Erdos–Renyi graph instantiation, the Erdos–Renyi using the Bernoulli parameter $\vartheta'G_n$ (the density of the connectome G_n). The resulting plots are in the right hand side of Fig. 7. Simple calculations of the distributions show that the pairs of graphs being seeded graph matched here in these repeated experiments are precisely correlated Erdos–Renyi graphs with the parameter ρ being precisely the edge correlation ϱ_e , which is equal to ϱ_T since $\varrho_h = 0$.

To emphasize: The left hand side of Fig. 7 is from seeded graph matching connectome to noisy connectome, and the right hand side of Fig. 7 is from seeded graph matching synthetic data of the same connectome density to a noisy version of this synthetic data, which turns out to precisely be seeded graph matching pairs of correlated Bernoulli random graphs where the noise parameter turns out to be the total correlation, so the figures in the right hand side of Fig. 7 are of Sect. “Of hockey sticks and broken hockey sticks” variety (except that the alignment strength values are averaged over 20 instantiations).

Notice that the figures in the left hand side of Fig. 7 and their respective counterparts in the right hand side of Fig. 7 look remarkably similar in many important ways. The differences seem to just be that the seeded graph matching success and alignment strength values clearly exhibit thresholding in the synthetic data, which is less pronounced and more gradual in the connectome data, although the sharpness of the connectome thresholding seems to be catching up as the number of vertices increases. Aside from this, there stills seems to be a reasonable phantom alignment strength for the connectome data.

Real data; matching same objects under different modalities

In this section, we illustrate the ideas in this paper using three real data sets from Fishkind et al. (2019a); they are the Wikipedia, Enron, and C Elegans pairs of graphs. Each is



an example of a pair of graphs with the same underlying objects (thus there is a natural “true” bijection), and adjacencies between objects in the respective graphs are relationships among the objects in two different modalities.

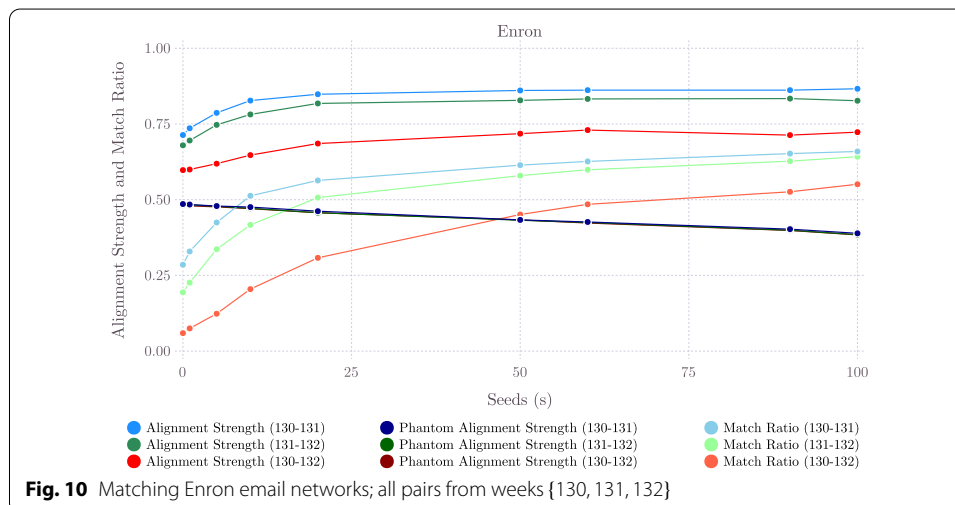
The Wikipedia pair of graphs G_1, G_2 from Fishkind et al. (2019a) were created in the year 2009. The vertices of G_1 are the English language Wikipedia articles hyperlinked from the Wikipedia article “Algebraic Geometry,” and all Wikipedia articles hyperlinked from these articles; in total, there are $n = 1382$ vertices. These vertices/articles each have directly corresponding articles in the French language Wikipedia, and these are the vertices of G_2 . Every pair of vertices/articles in G_1 are adjacent in G_1 precisely when one of the articles hyperlinks to the other article in the English language Wikipedia, and every pair of vertices in G_2 are adjacent in G_2 precisely when one of the articles links to the other in the French language Wikipedia. Thus G_1 and G_2 are simple, undirected graphs, and the “true” bijection is the function mapping English articles to their French versions.

For each value of $s = 0, 5, 50, 150, 250, 382, 500$, we did 100 replicates of uniformly sampling s seeds from the n vertices, seeded graph matched G_1 to G_2 using SGM, then recording the alignment strength $\text{str}(\hat{\varphi}_{SGM})$, averaged over the 100 replicates, plotted (in blue) vs the number of seeds s in Fig. 8. In the same figure, we recorded the match ratio (the number of nonseeds correctly matched, divided by the number of nonseeds), averaged over the 100 replications, plotted (in purple) versus the number of seeds s , also in Fig. 8. In addition, for each value of $s = 0, 5, 50, 150, 250, 382, 500$, we did 100 replicates of realizing uncorrelated pairs of Erdos–Renyi graphs H_1, H_2 , each with 1382 vertices and Bernoulli parameter of H_1 equal to the density of G_1 , Bernoulli parameter of H_2 equal to the density of G_2 , then uniformly sampling s seeds from the 1382 vertices, then seeded graph matched H_1 to H_2 using SGM, and recording the alignment strength $\text{str}(\hat{\varphi}_{SGM})$, averaged over the 100 replicates, plotted (in green) vs the number of seeds s in Fig. 8; these values represent the phantom alignment strength values in the respective seed levels. Note that, as the number seeds went from 0 to 5 to 50, the jump in match ratio coincides with a jump in the gap between seeded graph matching alignment strength and the phantom alignment strength. (Even when $s = 0$ there is some truth in the graph match; the match ratio was .0151, approximately 21 nonseed vertices matched correctly, whereas chance is $1/1382$, one nonseed vertex matched correctly.)

The *C. Elegans* pair of graphs G^{el} , G^{ch} from Varshney et al. (2011); Fishkind et al. (2019a) are connectomes mapping out the neural structure of the roundworm *Caenorhabditis Elegans*. *C. Elegans* is of interest to neuroscientists due to its well studied genetics (*C. Elegans* sequencing consortium 1998), comparatively simple nervous system (White et al. 1986), and a growing understanding of the correspondence between the two (Bargmann 1998; Arnatkevičiūtė et al. 2018). Like in humans, communication in the *C. Elegans* nervous system occurs via synapses, or junctions, between pairs of neurons. Neuronal synapses in the *C. Elegans* connectome can be classified in two ways (Varshney et al. 2011): an electrical synapse is a channel through which electrical impulses traverse, whereas chemical synapses are junctions through which neurotransmitters flow. We consider $n = 279$ somatic neurons of the hermaphrodite *C. Elegans* as the vertices of each graph. For each pair of vertices/neurons, they are adjacent in G^{el} precisely when there is an electrical synapse between them, and they are adjacent in G^{ch} precisely when there is a chemical synapse between them.

We conducted the identical experiments as we did for the Wikipedia graphs, except that the number of seeds s considered were $s = 0, 1, 5, 10, 20, 50, 75, 100, 150, 200$, and we matched G^{el} and G^{ch} . The resulting plots are in Fig. 9; alignment strength of the *C. Elegans* seeded graph match in blue, phantom alignment strength in green, match ratio in purple. Note that seeded graph matching did poorly, as evidenced by low match ratio, even when the number of seeds was huge (200 seeds and 79 nonseeds), and correspondingly the gap between seeded graph matching alignment strength and phantom alignment strength was small.

The Enron graphs from Fishkind et al. (2019a) arose in the following manner. Enron was a large and highly respected energy company that dissolved spectacularly in 2001 amid systemic fraud. The United States Justice Department released a trove of email messages between company employees. The graphs G_{130} , G_{131} , and G_{132} have as vertices $n = 184$ Enron employees and, for each pair of vertices/employees, the vertices are adjacent in G_{130} precisely when there is an email from one employee to the other in week number 130 of the email corpus, they are adjacent in G_{131} precisely when there is an email from one employee to the other in week number 131, and they are adjacent in G_{132}



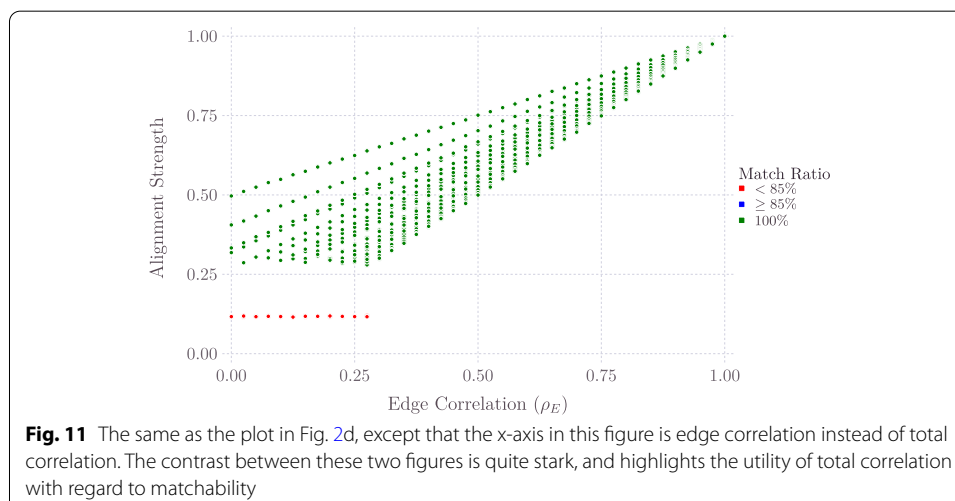
precisely when there is an email from one employee to the other in week number 132. The paper (Priebe et al. 2005) identified an anomaly going into week 132, and (Fishkind et al. 2019a) used match ratio differences between pairs of these graphs to highlight this anomaly.

We conducted the identical experiments for each of the pairs G_{130}, G_{131} and G_{131}, G_{132} and G_{130}, G_{132} as we did for the Wikipedia graphs, except that the number of seeds s considered were $s = 0, 1, 5, 10, 20, 50, 60, 90, 100$. The resulting plots are in Fig. 10. As noted in Fishkind et al. (2019a), the match ratio from matching G_{130} to G_{131} is highest of the three, since the anomaly had not yet occurred. The next highest match ratio was from matching G_{131} to G_{132} , then came matching G_{130} to G_{132} . Note that the gap between seeded graph matching alignment strength and phantom alignment strength was ordered the same way; highest was G_{130} to G_{131} , then was G_{131} to G_{132} , and then was G_{130} to G_{132} . Indeed, more gap here when there was higher match ratio. (Note that the match ratios here differ a bit from those in the paper (Fishkind et al. 2019a), Figure 8; that figure was inadvertently from a nonsimple graph version of the data, and here we created a simple graph.)

Notable mentions and future directions, plus caveats

The applications of graph matching are broad and many, and getting the right answer is only valuable when we know that we have the right answer. This paper provides principled tools that can help the practitioner decide if seeded graph matching has found the true bijection.

The first caveat—and future direction—is that we are presenting a conjecture, and not a theorem. Indeed, the Phantom Alignment Strength Conjecture, as formulated in Sect. “Phantom alignment strength conjecture, consequences,” includes terms in quotes; “moderate,” “high probability,” “very different,” and “negligible.” Ironing these terms out with specifics is part of the puzzle of proving the conjecture, and is an important next task. It may be a hard task, and we expect this paper to stimulate more experimentation, fine-tuning, and eventually a proof of the conjecture.



Part of the first caveat is the consideration that we don't have a proof of the conjecture as of now, and our experimentation is wide but not exhaustive, and thus there may be additional hypotheses or limitations to the conjecture statement.

A second caveat is that the conjecture is expressed in terms of an underlying model for a pair of random graphs, and we need to consider if particular real data that we may encounter (beyond the examples that we used here) can more generally be considered as arising from such a model.

Also, when there are multiple blocks with Bernoulli coefficients being realized from different distributions for different blocks, we saw in Sect. “[Block settings](#)” that total correlation became a much less reliable tool for determining matchability. More work is needed to explore this further; the paper (Fishkind et al. 2019b), when presenting empirical evidence for the relationship between total correlation and matchability, restricted their attention to the setting hypothesized in our Phantom Alignment Strength Conjecture, which excludes multiple blocks. Indeed, in the setting of our conjecture, the role of total correlation in matchability is starkly visible. See Fig. 2d, where the x-axis is total correlation, and compare to Fig. 11. Figure 11 is the same data plotted in Fig. 2d, except that the x-axis is used for edge correlation instead of total correlation. The contrast between these two figures is quite dramatic. Indeed, the current-literature-standard yardstick of edge correlation failed miserably in capturing matchability, whereas total correlation captured matchability perfectly here. These two figures are powerful illustration of the role of total correlation in matchability.

While more work remains to be done, we have here presented principled tools that can be of significant help to the practitioner now.

A large measure of inspiration for this paper came from Figures 1, 2, and 3 of Fishkind et al. (2019b) (on which half of us are co-authors). Those figures displayed the results of graph matchings of many simulations of pairs of correlated Bernoulli random graphs under similar conditions of the Phantom Alignment Strength Conjecture. One axis of each figure tracked edge correlation, and the second axis tracked heterogeneity correlation; green, yellow, and red dots were respectively located at coordinates corresponding to parameters where the graph matchings were always the truth, mostly the truth, and often not the truth, respectively. It was striking to observe that the regions of red and green were sharply demarcated by a level curve of total correlation, with little yellow between the red and green. These figures starkly demonstrated the role of total correlation in matchability, as well as thresholding behavior. Together with the theoretical results of Fishkind et al. (2019b) tying alignment strength to total correlation (when graph matching gets truth), we had important ingredients for the “hockey stick” at the heart of the Phantom Alignment Strength Conjecture.

Abbreviations

C. Elegans: *Caenorhabditis Elegans*; dMRI: Diffusion-weighted magnetic resonance imaging; DT-MRI: Diffusion tensor magnetic resonance imaging.

Acknowledgements

The authors thanks Dr. David Marchette for creating the Wikipedia graph that we used here.

Author's contributions

DEF formulated the Phantom Alignment Strength Conjecture, penned much of the manuscript. FP, HS, and LM did the numerical experiments and contributed to the formulation of the conjecture. VL penned the introduction section, including the literature review, and made the bibliography. EB was involved in the acquisition of the human connectome

data set, processed the data for our particular use, and penned a portion of the manuscript involving the human and *C. Elegans* connectome networks. VL, CEP, and AA worked on the theoretical development, edited the manuscript. All authors read and approved the manuscript.

Funding

This paper is based on research sponsored by the Air Force Research Laboratory and DARPA under Agreement Number FA8750-20-2-1001 and FA8750-17-2-0112. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the Air Force Research Laboratory and DARPA or the U.S. Government.

Availability of data and materials

Located at <https://cs.jhu.edu/~fparker9/phantom-alignment-strength/>.

Declarations

Competing interests

The authors declare that they have no competing interests.

Author details

¹Department of Applied Mathematics and Statistics, Johns Hopkins University, Baltimore, MD, USA. ²Center for Systems Science and Engineering, Johns Hopkins University, Baltimore, MD, USA. ³Department of Biostatistics, Johns Hopkins University, Baltimore, MD, USA. ⁴Department of Mathematics, University of Maryland, College Park, College Park, MD, USA.

Received: 8 March 2021 Accepted: 3 June 2021

Published online: 28 August 2021

References

- Arnatkevičiūtė A, Fulcher BD, Pocock R, Fornito A (2018) Hub connectivity, neuronal diversity, and gene expression in the *Caenorhabditis elegans* connectome. *PLoS Comput Biol* 14(2):1005989. <https://doi.org/10.1371/journal.pcbi.1005989>
- Arroyo J, Athreya A, Cape J, Chen G, Priebe CE, Vogelstein JT (2019) Inference for multiple heterogeneous networks with a common invariant subspace. *arXiv preprint arXiv:1906.10026*
- Babai L (2016) Graph isomorphism in quasipolynomial time. In: *Proceedings of the forty-eighth annual ACM symposium on theory of computing*. ACM, pp 684–697
- Barak B, Chou C, Lei Z, Schramm T, Sheng Y (2019) (nearly) efficient algorithms for the graph matching problem on correlated random graphs. In: *Advances in neural information processing systems*, pp 9190–9198
- Bargmann CI (1998) Neurobiology of the *Caenorhabditis elegans* genome. *Science* 282(5396):2028–2033. <https://doi.org/10.1126/science.282.5396.2028> (9851919)
- Bullmore E, Sporns O (2009) Complex brain networks: graph theoretical analysis of structural and functional systems. *Nat. Rev. Neurosci* 10:186–198
- Chen L, Vogelstein JT, Lyzinski V, Priebe CE (2016) A joint graph inference case study: the *C. elegans* chemical and electrical connectomes. In: *Worm*, vol 5. Taylor & Francis
- Chung J, Bridgeford E, Arroyo J, Pedigo BD, Saad-Eldin A, Gopalakrishnan V, Xiang L, Priebe CE, Vogelstein JT (2020) Statistical connectomics. *OSF Preprints*. <https://doi.org/10.31219/osf.io/ek4n3>
- Conte D, Foggia P, Sansone C, Vento M (2004) Thirty years of graph matching in pattern recognition. *Int J Pattern Recognit Artif Intell* 18(03):265–298
- Cullina D, Kiyavash N, Mittal P, Poor HV (2019) Partial recovery of Erdős–Rényi graph alignment via k-core alignment. *Proc ACM Meas Anal Comput Syst* 3(3):1–21
- Cullina D, Kiyavash N (2016) Improved achievability and converse bounds for Erdős–Rényi graph matching. In: *ACM SIGMETRICS performance evaluation review*, vol 44. ACM, pp 63–72
- Cullina D, Kiyavash N (2017) Exact alignment recovery for correlated Erdős–Rényi graphs. *CoRR* **abs/1711.06783**
- C. elegans* sequencing consortium: genome sequence of the nematode *C. elegans*: a platform for investigating biology. *Science* 282(5396):2012–2018 (1998). <https://doi.org/10.1126/science.282.5396.2012> 9851916
- Ding J, Ma Z, Wu Y, Xu J (2020) Efficient random graph matching via degree profiles. *Probab Theory Relat Fields*, pp 1–87
- Ding X, Zhang L, Wan Z, Gu M (2010) A brief survey on de-anonymization attacks in online social networks. In: *2010 International conference on computational aspects of social networks*. IEEE, pp 611–615
- Durante D, Dunson DB et al (2018) Bayesian inference and testing of group differences in brain networks. *Bayesian Anal* 13(1):29–58
- Emmert-Streib F, Dehmer M, Shi Y (2016) Fifty years of graph matching, network alignment and network comparison. *Inf Sci* 346–347:180–197
- Fan Z, Mao C, Wu Y, Xu J (2019) Spectral graph matching and regularized quadratic relaxations II: Erdős–Rényi graphs and universality. *arXiv preprint arXiv:1907.08883*
- Feizi S, Quon G, Recamonde-Mendoza M, Médard M, Kellis M, Jadbabaie A (2016) Spectral alignment of networks. *arXiv preprint arXiv:1602.04181*
- Fishkind DE, Adali S, Patsolic HG, Meng L, Singh D, Lyzinski V, Priebe CE (2019a) Seeded graph matching. *Pattern Recognit* 87:203–215
- Fishkind DE, Meng L, Sun A, Priebe CE, Lyzinski V (2019b) Alignment strength and correlation for graphs. *Pattern Recognit Lett* 125:295–302

- Foggia P, Percannella G, Vento M (2014) Graph matching and learning in pattern recognition in the last 10 years. *Int J Pattern Recognit Artif Intell* 28(01):1450001
- Gray WR, Bogovic JA, Vogelstein JT, Landman BA, Prince JL, Vogelstein RJ (2012) Magnetic resonance connectome automated pipeline: an overview. *Pulse IEEE* 3(2):42–48
- Heimann M, Shen H, Safavi T, Koutra D (2018) Regal: representation learning-based graph alignment. In: *Proceedings of the 27th ACM international conference on information and knowledge management*, pp 117–126
- Kiar G, Bridgeford EW, Gray Roncal WR, Chandrashekar V, Mhembere D, Ryman S, Zuo X-N, Margulies DS, Craddock RC, Priebe CE, Jung R, Calhoun VD, Caffo B, Burns R, Milham MP, Vogelstein JT (2018) A high-throughput pipeline identifies robust connectomes but troublesome variability. *bioRxiv*, 188706
- Kivelä M, Arenas A, Barthélemy M, Gleeson JP, Moreno Y, Porter MA (2014) Multilayer networks. *J Complex Netw* 2(3):203–271
- Levin K, Athreya A, Tang M, Lyzinski V, Park CEY (2017) Priebe: A central limit theorem for an omnibus embedding of multiple random graphs and implications for multiscale network inference. *arXiv preprint arXiv:1705.09355*
- Lin L, Liu X, Zhu S-C (2010) Layered graph matching with composite cluster sampling. *IEEE Trans Pattern Anal Mach Intell* 32(8):1426–1442
- Lyzinski V (2018) Information recovery in shuffled graphs via graph matching. *IEEE Trans Inf Theory* 64(5):3254–3273
- Lyzinski V, Sussman DL (2020) Matchability of heterogeneous networks pairs. *Inf Inference J IMA* 9(4):749–783
- Lyzinski V, Fishkind DE, Priebe CE (2014) Seeded graph matching for correlated Erdos–Renyi graphs. *J Mach Learn Res* 15:3513–3540
- Lyzinski V, Fishkind DE, Fiori M, Vogelstein JT, Priebe CE, Sapiro G (2016) Graph matching: relax at your own risk. *IEEE Trans Pattern Anal Mach Intell* 38(1):60–73
- Mhembere D, Roncal WG, Sussman D, Priebe CE, Burns R (2013) Computing scalable multivariate global invariants of large (brain-) graphs. In: *2013 IEEE global conference on signal and information processing, GlobalSIP 2013—proceedings*. <https://doi.org/10.1109/GlobalSIP.2013.6736874>
- Mossel E, Xu J (2020) Seeded graph matching via large neighborhood statistics. *Random Struct Algorithms* 57(3):570–611
- Onaran E, Garg S, Erkip E (2016) Optimal de-anonymization in random graphs with community structure. *arXiv preprint arXiv:1602.01409*
- Patsolic H, Adali S, Vogelstein JT, Park Y, Priebe CE, Li G, Lyzinski V (2014 (2019 major revision)) Seeded graph matching via joint optimization of fidelity and commensurability. *arXiv preprint arXiv:1401.3813*
- Pedarsani P, Grossglauser M (2011) On the privacy of anonymized networks. In: *Proceedings of the 17th ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, pp 1235–1243
- Priebe CE, Conroy JM, Marchette DJ, Park Y (2005) Scan statistics on enron graphs. *Comput Math Org Theory* 11(3):229–247
- Priebe CE, Park Y, Vogelstein JT, Conroy JM, Lyzinski V, Tang M, Athreya A, Cape J, Bridgeford E (2019) On a two-truths phenomenon in spectral graph clustering. *Proc Natl Acad Sci USA* 116(13):5995–6000. <https://doi.org/10.1073/pnas.1814462116>
- Shirani F, Garg S, Erkip E (2018) Matching graphs with community structure: a concentration of measure approach. In: *2018 56th annual allerton conference on communication, control, and computing (Allerton)*. IEEE, pp 1028–1035
- Singh R, Xu J, Berger B (2007) Pairwise global alignment of protein interaction networks by matching neighborhood topology. In: *Annual international conference on research in computational molecular biology*. Springer, pp 16–31
- Sussman DL, Park Y, Priebe CE, Lyzinski V (2019) Matched filters for noisy induced subgraph detection. *IEEE Trans Pattern Anal Mach Intell*
- Szklarczyk D, Franceschini A, Wyder S, Forslund K, Heller D, Huerta-Cepas J, Simonovic M, Roth A, Santos A, Tsafou KP et al (2015) STRING v10: protein-protein interaction networks, integrated over the tree of life. *Nucl Acids Res* 43(D1):447–452
- Tong AHY, Lesage G, Bader GD, Ding H, Xu H, Xin X, Young J, Berriz GF, Brost RL, Chang M et al (2004) Global mapping of the yeast genetic interaction network. *Science* 303(5659):808–813
- Umeyama S (1988) An eigendecomposition approach to weighted graph matching problems. *IEEE Trans Pattern Anal Mach Intell* 10(5):695–703
- Varshney LR, Chen BL, Paniagua E, Hall DH, Chklovskii DB (2011) Structural properties of the *Caenorhabditis elegans* neuronal network. *PLoS Comput Biol* 7(2):1001066. <https://doi.org/10.1371/journal.pcbi.1001066>
- Vogelstein JT, Conroy JM, Lyzinski V, Podrazik LJ, Kratzer SG, Harley ET, Fishkind DE, Vogelstein RJ, Priebe CE (2014) Fast approximate quadratic programming for graph matching. *PLoS ONE* 10(04)
- Vogelstein JT, Bridgeford EW, Pedigo BD, Chung J, Levin K, Mensh B, Priebe CE (2019) Connectal coding: discovering the structures linking cognitive phenotypes to individual histories. *Curr Opin Neurobiol* 55:199–212
- Wang Q, Mao Z, Wang B, Guo L (2017) Knowledge graph embedding: a survey of approaches and applications. *IEEE Trans Knowl Data Eng* 29(12):2724–2743
- Wasserman S, Faust K (1994) *Social network analysis: methods and applications*. Cambridge University Press, Cambridge
- White JG, Southgate E, Thomson JN, Brenner S (1986) The structure of the nervous system of the nematode *Caenorhabditis elegans*. *Philos Trans R Soc Lond Ser B* 314(1165):1–340. <https://doi.org/10.1098/rstb.1986.0056> (22462104)
- Yartseva L, Grossglauser M (2013) On the performance of percolation graph matching. In: *Proceedings of the first ACM conference on online social networks*, pp 119–130
- Zaslavskiy M, Bach F, Vert J-P (2009) A path following algorithm for the graph matching problem. *IEEE Trans Pattern Anal Mach Intell* 31(12):2227–2242
- Zhang S, Tong H (2016) Final: fast attributed network alignment. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp 1345–1354

- Zhou F, De la Torre F (2012) Factorized graph matching. In: 2012 IEEE conference on computer vision and pattern recognition (CVPR). IEEE, pp 127–134
- Zuo X-N, Anderson JS, Bellec P, Birn RM, Biswal BB, Blautzik J, Breitner JCS, Buckner RL, Calhoun VD, Castellanos FX, Chen A, Chen B, Chen J, Chen X, Colcombe SJ, Courtney W, Craddock RC, Di Martino A, Dong H-M, Fu X, Gong Q, Gorgolewski KJ, Han Y, He Y, He Y, He Y, Ho E, Holmes A, Hou X-H, Huckins J, Jiang T, Jiang Y, Kelley W, Kelly C, King M, LaConte SM, Lainhart JE, Lei X, Li H-J, Li K, Li K, Lin Q, Liu D, Liu J, Liu X, Liu Y, Lu G, Lu J, Luna B, Luo J, Lurie D, Mao Y, Margulies DS, Mayer AR, Meindl T, Meyerand ME, Nan W, Nielsen JA, O'Connor D, Paulsen D, Prabhakaran V, Qi Z, Qiu J, Shao C, Shehzad Z, Tang W, Villringer A, Wang H, Wang K, Wei D, Wei G-X, Weng X-C, Wu X, Xu T, Yang N, Yang Z, Zang Y-F, Zhang L, Zhang Q, Zhang Z, Zhang Z, Zhao K, Zhen Z, Zhou Y, Zhu X-T, Milham MP (2014) An open science resource for establishing reliability and reproducibility in functional connectomics. *Sci Data* 1(140049):1–13. <https://doi.org/10.1038/sdata.2014.49>

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Submit your manuscript to a SpringerOpen[®] journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)
