

482 **A Proofs and additional theoretical results**

483 **A.1 Proofs**

484 Here we usually omit the  $= \underline{k}$  suffix in  $\mathcal{A}_{q,t=\underline{k}}^{(s)}$  and write instead  $\mathcal{A}_{q,t}^{(s)}$  for easier reading. The operator  
 485  $\mathcal{A}_{q,t}^{(s)}$  is understood as having a fixed value  $t(x_s)$  of  $t$  observed. Similarly we often write  $\mathbb{E}_{q,t}$  to  
 486 abbreviate  $\mathbb{E}_{q,t=\underline{k}}$ , with the same understanding.

487 **Proof of Lemma 3.1** For convenience we repeat the statement of the lemma here.

488 **Lemma 3.1** *If  $q_{\underline{k}}(x) > 0$  for all  $\underline{k}$ , then  $\mathcal{A}_{q,t=\underline{k}}^{(s)}$  is a Stein operator for the conditional distribution of*  
 489  *$X$  given  $t(X) = \underline{k}$ , and  $\sum_s \mathcal{A}_{q,t=\underline{k}}^{(s)}$  is a Stein operator for the conditional distribution of  $X$  given*  
 490  *$t(X) = \underline{k}$ .*

491 *Proof.* In order to show the assertion we prove that for  $\mathbb{E}_{q,t=\underline{k}}$  denoting the conditional distribution of  
 492  $X$  given  $t(X) = \underline{k}$ , the expectation  $\mathbb{E}_{q,t=\underline{k}}[\mathcal{A}_{q,t}^{(s)}f]$  vanishes for all functions for which the expectation  
 493 exists. Let  $f$  be such a function. We have

$$\begin{aligned}\mathcal{A}_{q,t}^{(s)}f(x^{(s,1)}) &= q(x^{(s,0)}|t(x_s) = \underline{k})(f(x^{(s,0)}) - f(x^{(s,1)})) \\ \mathcal{A}_{q,t}^{(s)}f(x^{(s,0)}) &= q(x^{(s,1)}|t(x_s) = \underline{k})(f(x^{(s,1)}) - f(x^{(s,0)})).\end{aligned}$$

494 Thus,  $\mathbb{E}_{q,t}$ ,

$$\begin{aligned}\mathbb{E}_{q,t}[\mathcal{A}_{q,t}^{(s)}f] &= \sum_{x_{-s}} \left\{ \mathbb{1}(x_s = 1)p_t(x^{(s,1)})q(y^{(s,0)}|t(y_s) = \underline{k})(f(x^{(s,0)}) - f(x^{(s,1)})) \right. \\ &\quad \left. - \mathbb{1}(x_s = 0)p_t(x^{(s,0)})q(y^{(s,1)}|t(y_s) = \underline{k})(f(x^{(s,0)}) - f(x^{(s,1)})) \right\} \\ &= \sum_{x_{-s}} (f(x^{(s,0)}) - f(x^{(s,1)}))p_t(x^{(s,1)})p_t(x^{(s,0)}) \{ \mathbb{1}(x_s = 1) - \mathbb{1}(x_s = 0) \} \\ &= 0.\end{aligned}$$

495

□

496 **Proof of Theorem 3.2** For convenience we repeat the theorem here.

497 **Theorem 3.2** *Assume that  $\widehat{q}_t(x^{(s,1)})$  is a consistent estimator for  $q_t(x^{(s,1)})$  as  $L \rightarrow \infty$ . Then for any*  
 498 *function  $f$  such that  $\|\Delta f\| < \infty$  we have  $\mathbb{E}_q[\mathcal{A}_{\widehat{q},t}f(x)] \rightarrow \mathbb{E}_q[\mathcal{A}_{q,t}f(x)] = 0$  as  $L \rightarrow \infty$ .*

*Proof.* We recall the notation that Equation (4). We have that

$$\mathcal{A}_{q,t}f(x) = \frac{1}{N} \sum_{s \in [N]} [q(x^{(s,1)}|t(x_{-s}))f(x^{(s,1)}) + q(x^{(s,0)}|t(x_{-s}))f(x^{(s,0)}) - f(x)]$$

499 so that

$$\begin{aligned}\mathcal{A}_{\widehat{q},t}f(x) - \mathcal{A}_{q,t}f(x) &= \frac{1}{N} \sum_{s \in [N]} \{ (\widehat{q}(x^{(s,1)}|t(x_{-s})) - q(x^{(s,1)}|t(x_{-s})))f(x^{(s,1)}) \\ &\quad + (1 - \widehat{q}(x^{(s,1)}|t(x_{-s})) - (1 - q(x^{(s,1)}|t(x_{-s}))))f(x^{(s,0)}) \} \\ &= \frac{1}{N} \sum_{s \in [N]} (\widehat{q}(x^{(s,1)}|t(x_{-s})) - q(x^{(s,1)}|t(x_{-s})))\Delta_s f(x).\end{aligned}$$

Hence

$$|\mathcal{A}_{\widehat{q},t}f(x) - \mathcal{A}_{q,t}f(x)| \leq \|\Delta f\| \frac{1}{N} \sum_{s \in [N]} |\widehat{q}(x^{(s,1)}|t(x_{-s})) - q(x^{(s,1)}|t(x_{-s}))|.$$

500 Thus, if for all  $s \in [N]$ , as  $L \rightarrow \infty$  we have  $\widehat{q}(x^{(s,1)}|t(x_{-s})) - q(x^{(s,1)}|t(x_{-s})) \rightarrow 0$  for all  $s$  then so  
 501 does  $|\mathcal{A}_{\widehat{q},t}f(x) - \mathcal{A}_{q,t}f(x)|$ . The assertion follows from the assumption that  $\widehat{q}_t(x^{(s,1)})$  is a consistent  
 502 estimator for  $q_t(x^{(s,1)})$  as  $L \rightarrow \infty$ . □

503 **Proof of Theorem 3.3** For convenience we repeat the statement of the theorem here. Recall that a  
504 random graph model is *edge-exchangeable* if its edge indicator variables are finitely exchangeable.  
505 Often we just write *edge-exchangeable graph*. An ERGM is an example of an edge-exchangeable  
506 graph.

**Theorem 3.3** *If the graph is edge-exchangeable, then  $\text{AgraSS}t^2(\hat{q}, t; x)$  is a consistent estimator of*

$$\text{gKSS}^2(q; x) = N^{-2} \sum_{s, s' \in [N]} \left\langle \mathcal{A}_q^{(s)} K(x, \cdot), \mathcal{A}_q^{(s')} K(\cdot, x) \right\rangle_{\mathcal{H}}.$$

507 For easier tractability the proof is organised in two steps.

- 508 1. First, Proposition A.1 shows that in an edge-exchangeable random graph model,  $g_{\underline{k}}$  given in  
509 Equation (8) is a consistent estimator for  $q(x^{(s,1)} | t(x_{-s}))$  as  $NL \rightarrow \infty$ .
- 510 2. Theorem A.2 uses these results to obtain a concentration bound for  $\mathcal{A}_{\hat{q}}$  from which then  
511 Theorem 3.3 follows.

512 Moreover theoretical guarantees for fixed  $L$  which depend on the model are given. As the graph  
513 generator can generate as large a number  $L$  of graphs as desired, these theoretical results can be used  
514 to determine  $L$  which result in theoretical guarantees on deviations from the mean.

**Proposition A.1.** *Suppose that  $X_1, \dots, X_L$  are i.i.d. copies of the adjacency matrix of an edge-exchangeable random graph model. Let  $s = (i, j)$  be a fixed vertex-pair. For  $l = 1, \dots, L$  and for a graph  $X_l$  let  $t_l^{(s)}$  denote a possibly multivariate statistic which is evaluated on the collection of indicator variables in  $X_l$  except  $X_{s,l}$ . For a possible  $t_l^{(s)}$  outcome  $\underline{k}$ , let  $p(\underline{k}) = \mathbb{P}(t_l^{(s)} = \underline{k})$  and let  $\underline{k}$  be such that  $p(\underline{k}) \neq 0$ . Set*

$$p(1; \underline{k}) = \mathbb{P}(X_s = 1 | t^{(s)} = \underline{k});$$

let

$$n(\underline{k}, s) = \sum_{l=1}^L X_{s,l} \mathbb{1}(t_l^{(s)} = \underline{k}) \quad \text{and} \quad n(\underline{k}^{(s)}) = \sum_{l=1}^L \mathbb{1}(t_l^{(s)} = \underline{k});$$

$$n(\underline{k}) = \sum_{s \in [N]} n(\underline{k}, s) \quad \text{and} \quad N_{\underline{k}} = \sum_{s \in [N]} n(\underline{k}^{(s)});$$

and set

$$g(\underline{k}) = \frac{n(\underline{k})}{N_{\underline{k}}} \mathbb{1}(N_{\underline{k}} \geq 1).$$

Then  $g(\underline{k}) \rightarrow p$  in probability as  $NL \rightarrow \infty$ . In particular, for all  $\epsilon > 0$ ,

$$\mathbb{P} [|\hat{g}(\underline{k}) - p(1; \underline{k})| > \epsilon] \leq \frac{4}{\epsilon^2 NL} \{p(s, \underline{k})[1 - p(s, \underline{k})p(\underline{k})] + 1 - p(\underline{k})\}.$$

*Proof.* Due to the exchangeability of the edges we have, with  $s$  denoting a generic edge,

$$\mathbb{E}(n(\underline{k})) = N\mathbb{E}(n(\underline{k}, s)) = NLp(1; \underline{k})p(\underline{k}), \quad \text{Var}(n(\underline{k})) = NLp(1; \underline{k})p(\underline{k})(1 - p(1; \underline{k})p(\underline{k})),$$

as well as

$$\mathbb{E}(N_{\underline{k}}) = N\mathbb{E}(n(\underline{k}^{(s)})) = NLp(\underline{k}), \quad \text{Var}(N_{\underline{k}}) = NLp(\underline{k})(1 - p(\underline{k})).$$

515 To show convergence in probability, let  $\epsilon > 0$ . Then

$$\mathbb{P} [|\hat{g}(\underline{k}) - p(1, \underline{k})| > \epsilon] \leq \mathbb{P} \left[ \left| \frac{n(\underline{k})}{\mathbb{E}(n(\underline{k}))} \mathbb{1}[n(\underline{k}) \geq 1] - p(1, \underline{k}) \right| > \frac{1}{2} \epsilon \right]$$

$$+ \mathbb{P} \left[ \left| n(\underline{k}) \mathbb{1}[n(\underline{k}) \geq 1] \left( \frac{1}{n(\underline{k})} - \frac{1}{\mathbb{E}(n(\underline{k}))} \right) \right| > \frac{1}{2} \epsilon \right].$$

516 Note that  $n(\underline{k})\mathbb{1}[n(\underline{k}) \geq 1] = n(\underline{k})$ . By Chebychev's inequality,

$$\begin{aligned} \mathbb{P} \left[ \left| \frac{n(\underline{k})}{\mathbb{E}(N_{\underline{k}})} \mathbb{1}[N_{\underline{k}} \geq 1] - p(1; \underline{k}) \right| > \frac{1}{2} \epsilon \right] &\leq \frac{4}{\epsilon^2 N^2 L^2 p(\underline{k})^2} N L p(1; \underline{k}) p(\underline{k}) (1 - p(1; \underline{k}) p(\underline{k})) \\ &= \frac{4}{\epsilon^2 N L p(\underline{k})} p(1; \underline{k}) (1 - p(1; \underline{k}) p(\underline{k})) \end{aligned}$$

517 and

$$\begin{aligned} \mathbb{P} \left[ \left| n(\underline{k}) \mathbb{1}[N_{\underline{k}} \geq 1] \left( \frac{1}{n(\underline{k})} - \frac{1}{\mathbb{E}(N_{\underline{k}})} \right) \right| > \frac{1}{2} \epsilon \right] &\leq \mathbb{P} \left[ \frac{1}{\mathbb{E}(N_{\underline{k}})} |\mathbb{E}(N_{\underline{k}}) - N_{\underline{k}}| > \frac{1}{2} \epsilon \right] \\ &\leq \frac{4}{\epsilon^2 N L p(\underline{k})} (1 - p(\underline{k})). \end{aligned}$$

518 Summing the contributions completes the proof.  $\square$

Proposition A.1 shows that in edge-exchangeable graphs,  $\hat{g}(\underline{k})$  consistently estimates  $q(x^{(s,1)}) = q(x^{(s,1)} | t(x_{-s}) = \underline{k})$ . In an expanded version of Theorem 3.3 we show that the approximate Stein operator from Eq.(7),

$$\mathcal{A}_{\hat{q},t} f(x) := \frac{1}{N} \sum_{s \in [N]} \mathcal{A}_{\hat{q}(x^{(s)} | t(x_{-s}))} f(x),$$

with

$$\hat{q}(x^{(s,1)} | t(x_{-s})) = g(x_{-s}), \quad \hat{q}(x^{(s,0)} | t(x_{-s})) = 1 - g(x_{-s})$$

is a consistent estimator of

$$\mathcal{A}_{q,t} f(x) := \frac{1}{N} \sum_{s \in [N]} \mathcal{A}_{q(x^{(s)} | t(x_{-s}))} f(x).$$

We recall

$$\text{AgraSS}^2(\hat{q}, t, x) = \frac{1}{N^2} \sum_{s, s' \in [N]} h_x(s, s')$$

with

$$h_x(s, s') = \left\langle \mathcal{A}_{\hat{q},t}^{(s)} K(x, \cdot), \mathcal{A}_{\hat{q},t}^{(s')} K(\cdot, x) \right\rangle_{\mathcal{H}}.$$

519 We state the expanded version of Theorem 3.3 here.

**Theorem A.2.** *If the graph is edge-exchangeable then for any test function  $f$  for which the Stein operator  $\mathcal{A}_{q,t} f$  is well defined, and for all  $\epsilon > 0$*

$$\mathbb{P}(|\mathcal{A}_{\hat{q},t} f(X) - \mathcal{A}_{q,t} f(X)| > \epsilon) \leq \frac{4}{\epsilon^2 N L (|\Delta f|)^{-2}} (\{p(1, \underline{k})[1 - p(1, \underline{k})p(\underline{k})] + 1 - p(\underline{k})\}).$$

Moreover,  $\text{AgraSS}^2(\hat{q}, t, x)$  is a consistent estimator of

$$\text{gKSS}(x) = \frac{1}{N^2} \sum_{s, s' \in [N]} \left\langle \mathcal{A}_{q,t}^{(s)} K(x, \cdot), \mathcal{A}_{q,t}^{(s')} K(\cdot, x) \right\rangle_{\mathcal{H}}.$$

*Proof.* We have that

$$\mathcal{A}_{q,t} f(x) := \frac{1}{N} \sum_{s \in [N]} \mathcal{A}_{q(x^{(s)} | t(x_{-s}))} f(x).$$

and

$$\mathcal{A}_{\hat{q},t} f(x) := \frac{1}{N} \sum_{s \in [N]} \mathcal{A}_{\hat{q}(x^{(s)} | t(x_{-s}))} f(x),$$

with

$$\hat{q}(x^{(s,1)} | t(x_{-s})) = g(x_{-s}), \quad \hat{q}(x^{(s,0)} | t(x_{-s})) = 1 - g(x_{-s})$$

520 so that

$$\begin{aligned}
\mathcal{A}_{\hat{q},t}f(x) - \mathcal{A}_{q,t}f(x) &= \frac{1}{N} \sum_{s \in [N]} \{ (g(\underline{k}) - q(x^{(s)}|t(x_{-s})))f(x^{(s,1)}) \\
&\quad + (1 - g(\underline{k}) - (1 - q(x^{(s)}|t(x_{-s}))))f(x^{(s,0)}) \} \\
&= \frac{1}{N} \sum_{s \in [N]} (g(\underline{k}) - q(x^{(s)}|t(x_{-s}))) \{ f(x^{(s,1)}) - f(x^{(s,0)}) \} \\
&= \frac{1}{N} \sum_{s \in [N]} (g(\underline{k}) - q(x^{(s)}|t(x_{-s}))) \Delta_s f(x).
\end{aligned}$$

Hence

$$|\mathcal{A}_{\hat{q},t}f(x) - \mathcal{A}_{q,t}f(x)| \leq \|\Delta f\| \frac{1}{N} \sum_{s \in [N]} |g(\underline{k}) - q(x^{(s)}|t(x_{-s}))|.$$

521 With Proposition A.1 and using the edge-exchangeability,

$$\mathbb{P}(|\mathcal{A}_{\hat{q},t}f(X) - \mathcal{A}_{q,t}f(X)| > \epsilon) \leq \frac{4}{\epsilon^2 N L (|\Delta f|)^{-2}} (\{p(1, \underline{k})[1 - p(1, \underline{k})p(\underline{k})] + 1 - p(\underline{k})\}).$$

The fact that taking the sup over functions in the Hilbert space  $\mathcal{H}$  does not spoil the convergence follows from the closed form representation of the sup of AgraSSt<sup>2</sup>, see for example Equation (11) in [Xu and Reinert, 2021]. We have that

$$\text{AgraSSt}^2(\hat{q}, t, x) = \frac{1}{N^2} \sum_{s, s' \in [N]} h_x(s, s')$$

where

$$h_x(s, s') = \left\langle \mathcal{A}_{\hat{q},t}^{(s)} K(x, \cdot), \mathcal{A}_{\hat{q},t}^{(s')} K(\cdot, x) \right\rangle_{\mathcal{H}}.$$

522 Hence,

$$\begin{aligned}
&\text{AgraSSt}^2(\hat{q}) - \text{AgraSSt}^2(q) \\
&= \frac{1}{N^2} \sum_{s, s' \in [N]} \left\langle \mathcal{A}_{\hat{q},t}^{(s)} K(x, \cdot) - \mathcal{A}_{q,t}^{(s)} K(x, \cdot), \mathcal{A}_{\hat{q},t}^{(s')} K(\cdot, x) - \mathcal{A}_{q,t}^{(s')} K(\cdot, x) \right\rangle_{\mathcal{H}}
\end{aligned}$$

523 and the first part gives the desired convergence as  $L \rightarrow \infty$ . □

## 524 A.2 Gaussian approximation for AgraSSt in ERGMs

525 As ERGMs are edge-exchangeable models, Theorem A.2 shows that the AgraSSt operator is a  
526 consistent estimator for the ERGM Glauber Stein operator. If the observed graph  $x$  is a realisation of  
527 an ERGM then results from Xu and Reinert [2021] can be leveraged to obtain finer theoretical results.

First we detail the scaling for exponential random graph models which is used in the theoretical results which follow. For a graph  $H$  on at most  $n$  vertices  $V(H)$  denote the vertex set, and for  $x \in \{0, 1\}^N$ , denote by  $t(H, x)$  the number of *edge-preserving* injections from  $V(H)$  to  $V(x)$ ; an injection  $\sigma$  preserves edges if for all edges  $vw$  of  $H$  with  $\sigma(v) < \sigma(w)$ ,  $x_{\sigma(v)\sigma(w)} = 1$ . For  $v_H = |V(H)| \geq 3$  set

$$t_H(x) = \frac{t(H, x)}{n(n-1) \cdots (n - v_H + 3)}.$$

528 If  $H = H_1$  is a single edge, then  $t_H(x)$  is twice the number of edges of  $x$ . In the exponent this  
529 scaling of counts matches Definition 1 in Bhamidi et al. [2011] and Sections 3 and 4 of Chatterjee  
530 and Diaconis [2013]. An ERGM for the collection  $x \in \{0, 1\}^N$  can be defined as follows.

**Definition A.3** (Definition 1.5 in Reinert and Ross [2019]). Fix  $n \in \mathbb{N}$  and  $k \in \mathbb{N}$ . Let  $H_1$  be a single edge and for  $l = 2, \dots, k$  let  $H_l$  be a connected graph on at most  $n$  vertices; set  $t_l(x) = t_{H_l}(x)$ . For

$\beta = (\beta_1, \dots, \beta_k)^\top \in \mathbb{R}^k$  and  $t(x) = (t_1(x), \dots, t_k(x))^\top \in \mathbb{R}^k$ .  $X \in \mathcal{G}_n^{lab}$  follows the exponential random graph model  $X \sim \text{ERGM}(\beta, t)$  if for  $\forall x \in \mathcal{G}_n^{lab}$ ,

$$q(X = x) = \frac{1}{\kappa_n(\beta)} \exp\left(\sum_{l=1}^k \beta_l t_l(x)\right).$$

531 Here  $\kappa_n(\beta)$  is a normalisation constant.

In particular, under suitable conditions, the ERGM Glauber Stein operator is close to the  $G(n, p)$  Stein operator. This result is already shown in Reinert and Ross [2019], Theorem 1.7, with details provided in the proof of Theorem 1 in Xu and Reinert [2021]. To give the result, a technical assumption is required, which originates in Chatterjee and Diaconis [2013], and is required in Reinert and Ross [2019]. For  $a \in [0, 1]$ , define the following functions [Bhamidi et al., 2011, Eldan and Gross, 2018], with the notation in Definition A.3 for  $\text{ERGM}(\beta, t)$ :

$$\Phi(a) := \sum_{l=1}^k \beta_l e_l a^{e_l-1}, \quad \varphi(a) := \frac{1 + \tanh(\Phi(a))}{2}$$

532 where  $e_l$  is the number of edges in  $H_l$ .

533 *Assumption 1.* (1)  $\frac{1}{2} |\Phi'(1)| < 1$ . (2)  $\exists a^* \in [0, 1]$  that solves the equation  $\varphi(a^*) = a^*$ .

534 The value  $a^*$  will be the edge probability in the approximating Bernoulli random graph,  $\text{ER}(a^*)$ .

535 The following result holds.

**Proposition A.4.** *Let  $q(x) = \text{ERGM}(\beta, t)$  satisfy Assumption 1 and let  $\tilde{q}$  denote the distribution of  $\text{ER}(a^*)$ . Then there is an explicit constant  $C = C(\beta, t, K)$  such that for all  $\epsilon > 0$ ,*

$$\frac{1}{N} \sum_{s \in N} \mathbb{E} |(\mathcal{A}_q^{(s)} f(Y) - \mathcal{A}_{\tilde{q}}^{(s)} f(Y))| \leq \|\Delta f\| \binom{n}{2} \frac{C(\beta, t)}{\sqrt{n}}.$$

536 *Moreover, for  $f \in \mathcal{H}$  equipped with kernel  $K$ , let  $f_x^*(\cdot) = \frac{(\mathcal{A}_q - \mathcal{A}_{\tilde{q}})K(x, \cdot)}{\|(\mathcal{A}_q - \mathcal{A}_{\tilde{q}})K(x, \cdot)\|_{\mathcal{H}}}$ . Then there is an*  
 537 *explicit constant  $C = C(\beta, t, K)$  such that for all  $\epsilon > 0$ ,*

$$\mathbb{P}(|\text{gKSS}(q, X) - \text{gKSS}(\tilde{q}, Y)| > \epsilon) \tag{12}$$

$$\leq \left\{ \|\Delta(\text{gKSS}(q, \cdot))^2\| (1 + \|\Delta \text{gKSS}(q, \cdot)\|) + 4 \sup_x (\|\Delta f_x^*\|^2) \right\} \binom{n}{2} \frac{C}{\epsilon^2 \sqrt{n}}. \tag{13}$$

538 *Proof.* The assertion follows immediately from the proof of Theorem 1 in Xu and Reinert [2021].  $\square$

539 The approximation with a Bernoulli random graph is useful as for a Bernoulli random graphs a  
 540 normal approximation for its gKSS is available in Xu and Reinert [2021], under suitable assumptions.

541 *Assumption 2.* Let  $\mathcal{H}$  be the RKHS associated with the kernel  $K : \{0, 1\}^N \times \{0, 1\}^N \rightarrow \mathbb{R}$  and for  
 542  $s \in [N]$  let  $\mathcal{H}_s$  be the RKHS associated with the kernel  $l_s : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}$ . Then

543 i)  $\mathcal{H}$  is a tensor product RKHS,  $\mathcal{H} = \otimes_{s \in [n]} \mathcal{H}_s$ ;

544 ii)  $k$  is a product kernel,  $k(x, y) = \otimes_{s \in [N]} l_s(x_s, y_s)$ ;

545 iii)  $\langle l_s(x_s, \cdot), l_s(x_s, \cdot) \rangle_{\mathcal{H}_s} = 1$ ;

546 iv)  $l_s(1, \cdot) - l_s(0, \cdot) \neq 0$  for all  $s \in [N]$ .

547 These assumptions are satisfied for example for the suitably standardised Gaussian kernel  $K(x, y) =$   
 548  $\exp\{-\frac{1}{\sigma^2} \sum_{s \in [N]} (x_s - y_s)^2\}$ .

549 Letting  $\|\cdot\|_1$  denote  $L_1$ -distance, and  $\mathcal{L}$  denote the law of a random variable, Xu and Reinert [2021]  
 550 show the following normal approximation.

551 **Theorem A.5** (Theorem 2 in Xu and Reinert [2021]). *Let  $Y$  have the distribution  $\tilde{q}$  of a Bernoulli*  
 552 *random graph  $\text{ER}(a^*)$  as in Proposition A.4. Assume that the conditions i) - iv) in Assumption 2*  
 553 *hold. Let  $\mu = \mathbb{E}[\text{gKSS}^2(\tilde{q}, Y)]$  and  $\sigma^2 = \text{Var}[\text{gKSS}^2(\tilde{q}, Y)]$ . Set  $W = \frac{1}{\sigma}(\text{gKSS}^2(\tilde{q}, Y) - \mu)$  and*

554 let  $Z$  denote a standard normal variable, Then there is an explicit constant  $C = C(a^*, l_s, s \in [N])$   
 555 such that

$$\|\mathcal{L}(W) - \mathcal{L}(Z)\|_1 \leq \frac{C}{\sqrt{N}}.$$

556 Thus a normal approximation for the approximating gKSS can then be used to assess the theoretical  
 557 behaviour of AgraSSt as follows.

558 **Corollary A.6.** *Let the assumptions Proposition A.4 and Theorem A.5 be satisfied. With the notation*  
 559 *of Theorem A.5, assume that the RKHS kernel  $K$  is such that the right hand side of Proposition A.4*  
 560 *is  $o(n)$ . Then  $\frac{1}{\sigma}(\text{AgraSSt}(\tilde{q}(x^{(s)}|t(x_{-s}))) - \mu)$  is approximately standard normally distributed as*  
 561  *$N \rightarrow \infty$ .*

562 *Proof.* For all  $\epsilon > 0$ ,

$$\begin{aligned} & \mathbb{P} \left[ \left| \text{AgraSSt}(\tilde{q}(x^{(s)}|t(x_{-s}))) - \text{gKSS}(a^*) \right| > \epsilon \right] \\ & \leq \mathbb{P} \left[ \left| \text{AgraSSt}(\tilde{q}(x^{(s)}|t(x_{-s}))) - \text{gKSS}(q) \right| > \frac{1}{2}\epsilon \right] + \mathbb{P} \left[ \left| \text{gKSS}(q) - \text{gKSS}(a^*) \right| > \frac{1}{2}\epsilon \right]. \end{aligned}$$

563 The first summand tends to 0 as  $N \rightarrow \infty$  due to Theorem 3.2 and the second summand tends to 0  
 564 due to Proposition A.4. That  $\text{gKSS}(a^*)$  is approximately normally distributed with the appropriate  
 565 scaling follows from Theorem A.5.  $\square$

566 The theoretical behaviour of the subsampling version  $\widehat{\text{AgraSSt}}(\tilde{q}(x^{(s)}|t(x_{-s})))$  is addressed in  
 567 Proposition 3.4. A detailed examination of the choice of kernel  $K$  such that the assumptions of  
 568 Corollary A.6 are satisfied is left for future work.

## 569 B Additional background

570 In this section, we present additional background to complement the discussions in the main text.

### 571 B.1 Parameter estimation for random graphs

572 Estimating parameters for parametric models is possible *only* when the parametric family is *explicitly*  
 573 *specified*. For instance, in the synthetic example for E2ST model shown in Section 5.1,  $\hat{\beta}_l$  can  
 574 be estimated for  $\beta_l$  since the edge, 2Star and triangle statistics are specified. There are various  
 575 approaches for parameter estimation.

576 **Maximum likelihood** Maximum likelihood is a popular approach for parameter estimation in  
 577 random graph models. A complication arises because its probability mass function from Eq.(1),

$$q(X = x) = \frac{1}{\kappa_n(\beta)} \exp \left( \sum_{l=1}^k \beta_l t_l(x) \right).$$

578 involves a normalisation constant  $\kappa_n(\beta) = \sum_x \exp \{ \sum_{l=1}^k \beta_l t_l(x) \}$  which is generally intractable  
 579 and needs to be estimated for performing MLE. For this task, Markov chain Monte-Carlo maximum  
 580 likelihood estimation (MCMCMLE) for ERGM has been developed by Snijders [2002]. When  
 581 the network size is large, accurate estimation for the normalised  $\kappa_n(\beta)$  requires large amount of  
 582 Monte-Carlo samples and is hence computationally expensive.

583 **Maximum pseudo-likelihood estimator** To alleviate the problem associated with the normalising  
 584 constant, Maximum Pseudo-likelihood Estimation (MPLE) [Besag, 1975] has been developed for  
 585 ERGMs, see Strauss and Ikeda [1990] and also Schmid and Desmarais [2017]. MPLE factorises the  
 586 conditional edge probability to approximate the exact likelihood,

$$q(x) = \prod_{s \in [N]} q(x^s | x_{-s}). \quad (14)$$

587 For ERGMs the conditional distribution  $q(x^s|x_{-s})$  does not involve the normalising constant and  
 588 can hence be computed more efficiently than the MLE. However, in general the MPLS is not  
 589 consistent for ERGMs as the edges are generally non-independent. The consistency of MPLE for  
 590 Boltzmann machines is shown in Hyvärinen [2006]. A thorough comparison of MCMCMLE and  
 591 MPLE estimation in ERGMs can be found in Van Duijn et al. [2009].

592 **Contrastive divergence** Estimation based on contrastive divergence (CD) [Hinton, 2002] has also  
 593 been developed for ERGM estimation [Hunter and Handcock, 2006]. Contrastive divergence runs a  
 594 small number of Markov chains simultaneously for  $T$  steps and estimates the gradient based on the  
 595 differences between initial values and values after  $T$  steps in order to find a maximum. Convergence  
 596 results for exponential family models are shown in Jiang et al. [2018]. CD can provide a useful  
 597 balance between computationally expensive but accurate MCMCMLE and fast but inconsistent  
 598 MPLE.

## 599 B.2 Kernel Stein discrepancies and kernel-based nonparametric hypothesis testing

600 The task of hypothesis testing involves the comparison of distributions  $p$  and  $q$  that are significantly  
 601 different with respect to the size of the test, denoted by  $\alpha$ . In nonparametric tests, the distributions  
 602 are not assumed to be in any parametric families and test statistics are often based on ranking of  
 603 observations. In contrast, parametric tests, such as a Student t-test or a normality test, assume a  
 604 pre-defined parametric family to be tested against and usually employ a particular summary statistics  
 605 such as means or standard deviations. Recent advances in nonparametric test procedures introduce  
 606 RKHS functions which can be rich enough to distinguish distributions whenever they differ. Below  
 607 we detail two instances which are relevant for the main paper.

608 We start with a terse review of **kernel Stein discrepancy (KSD) for continuous distributions**  
 609 developed to compare and test distributions [Gorham and Mackey, 2015, Ley et al., 2017]. Let  $q$   
 610 be a smooth probability density on  $\mathbb{R}^d$  that vanishes at the boundary. The operator  $\mathcal{A}_q : (\mathbb{R}^d \rightarrow$   
 611  $\mathbb{R}^d) \rightarrow (\mathbb{R}^d \rightarrow \mathbb{R})$  is called a *Stein operator* if the following *Stein identity* holds:  $\mathbb{E}_q[\mathcal{A}_q f] = 0$ ,  
 612 where  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is any bounded smooth function. A suitable function class  $\mathcal{F}$  is such that if  
 613  $\mathbb{E}_p[\mathcal{A}_q f] = 0$  for all functions  $f \in \mathcal{F}$ , then  $p = q$  follows. It is convenient to take  $\mathcal{F} = B_1(\mathcal{H})$ , the  
 614 unit ball of a large enough RKHS with bounded kernel  $K$ . The kernel Stein discrepancy (KSD)  
 615 between two densities  $p$  and  $q$  based on  $\mathcal{A}_q$  is defined as

$$\text{KSD}(p||q, \mathcal{H}) = \sup_{f \in B_1(\mathcal{H})} \mathbb{E}_p[\mathcal{A}_q f]. \quad (15)$$

616 Under mild regularity conditions, for a particular choice of  $\mathcal{A}$  called Langevin operator,  
 617  $\text{KSD}(p||q, \mathcal{H}) \geq 0$  and  $\text{KSD}(p||q, \mathcal{H}) = 0$  if and only if  $p = q$  [Chwialkowski et al., 2016], in  
 618 which case KSD is a proper discrepancy measure between probability densities.

619 The KSD in Eq.(15) can be used to test the model goodness-of-fit as follows. One can show that  
 620  $\text{KSD}^2(p||q, \mathcal{H}) = \mathbb{E}_{x, \tilde{x} \sim p}[h_q(x, \tilde{x})]$ , where  $x$  and  $\tilde{x}$  are independent random variables with density  $p$   
 621 and  $h_q(x, \tilde{x})$  is given in explicit form which does not involve  $p$ ,

$$h_q(x, \tilde{x}) = \langle \mathcal{A}_q K(x, \cdot), \mathcal{A}_q K(\cdot, \tilde{x}) \rangle_{\mathcal{H}}. \quad (16)$$

622 Given a set of samples  $\{x_1, \dots, x_n\}$  from an unknown density  $p$  on  $\mathbb{R}^d$ , to test whether  $p = q$ , the  
 623 statistic  $\text{KSD}^2(p||q, \mathcal{H})$  can be empirically estimated by independent samples from  $p$  using a  $U$ - or  $V$ -  
 624 statistic. The critical value is determined by bootstrap based on weighted chi-square approximations  
 625 for  $U$ - or  $V$ -statistics. For goodness-of-fit tests of discrete distributions when i.i.d. samples are  
 626 available, a kernel discrete Stein discrepancy (KDSD) has been proposed in Yang et al. [2018].

**Goodness-of-fit Testing** aims to check the null hypothesis  $\mathcal{H}_0 : p = q$  against the general alternative  
 $\mathcal{H}_1 : p \neq q$  when the target distribution  $q$  is explicitly specified. Given sample(s) from the *unknown*  
 distribution  $p$  and an explicit density  $q$ ,  $\mathcal{H}_0$  is assessed using a chosen test statistic, usually a  
 discrepancy measure,  $D(q||p)$ , between  $p$  and  $q$ , which can be estimated empirically. Kernel-based  
 hypothesis tests on goodness-of-fit for continuous distributions  $q$  use the kernel Stein discrepancy

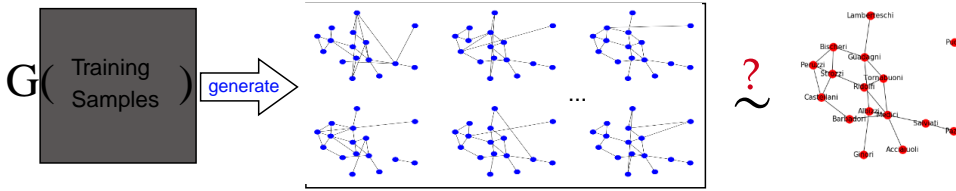


Figure 3: Assessing trained graph generators.

(KSD) in Section B.2 as the test statistic. Given samples  $x_1, \dots, x_n$  from the *unknown* density  $p$ ,  $\text{KSD}^2(p||q, \mathcal{H})$  in Eq.(15) is estimated via the  $V$ -statistic

$$\widehat{\text{KSD}}^2(p||q, \mathcal{H}) = \frac{1}{n^2} \sum_{i,j} h_q(x_i, x_j);$$

627 recall that  $h_q(x_i, x_j) = \langle \mathcal{A}_q K(x_i, \cdot), \mathcal{A}_q K(x_j, \cdot) \rangle_{\mathcal{H}}$  from Eq.(16). The null distribution of this test  
 628 statistic involves integral operators that are not available in close form; often it is simulated using  
 629 a wild-bootstrap procedure [Chwialkowski et al., 2014]. With the (simulated) null distribution, the  
 630 critical value of the test can be estimated to decide whether the null hypothesis is rejected at test  
 631 level  $\alpha$ . In this way, a general method for nonparametric testing of goodness-of-fit on  $\mathbb{R}^d$  is obtained,  
 632 which is applicable even for models with an intractable normalising constant.

633 **Two-sample Testing** aims to determine whether two sets of samples are drawn from the same  
 634 distribution, i.e. instead of  $q$  being available in density form as in the goodness-of-fit setting,  $q$   
 635 is only accessible through samples. Maximum mean embedding (MMD) test are often used for this  
 636 two-sample problem [Gretton et al., 2007]. These tests are based on the kernel mean embedding of a  
 637 distribution,

$$\mu_p := \mathbb{E}_{x \sim p}[k(x, \cdot)] = \int_{\mathcal{X}} k(x, \cdot) dp(x) \in \mathcal{H}, \quad (17)$$

638 whenever  $\mu_p$  exist. Similar to KSD, MMD takes the supremum over unit ball RKHS functions;

$$\text{MMD}(p||q) = \sup_{f \in B_1(\mathcal{H})} |\mathbb{E}_p[f] - \mathbb{E}_q[f]| = \|\mu_p - \mu_Q\|_{\mathcal{H}}. \quad (18)$$

639 With samples  $x_1, \dots, x_m \sim p$  and  $y_1, \dots, y_n \sim q$ , MMD can be estimated empirically via  $U$ -statistics,

$$\widehat{\text{MMD}}_u^2(p||q) = \frac{1}{m(m-1)} \sum_{i \neq i'} k(x_i, x_{i'}) + \frac{1}{n(n-1)} \sum_{j \neq j'} k(y_j, y_{j'}) - \frac{2}{mn} \sum_{ij} k(x_i, y_j). \quad (19)$$

640 In such kernel-based two-sample tests, the null distribution can be obtained via a permutation  
 641 procedure [Gretton et al., 2007]; this procedure can be more robust compared to a wild-bootstrap  
 642 procedure, especially when the kernels need to be optimised [Gretton et al., 2012, Jitkrittum et al.,  
 643 2016, Liu et al., 2020, 2021].

644 The two-sample procedure can also be applied to verify model assumptions when the model is not  
 645 directly accessible through its distribution but through generated samples. Such a strategy has been  
 646 considered as benchmark testing procedure in various studies for goodness-of-fit tests [Jitkrittum et al.,  
 647 2017, Xu and Matsuda, 2020, 2021]. Despite lower test power compared to the corresponding state-of-  
 648 the-art KSD-based tests and higher computational cost due to additional empirical estimation for the  
 649 distribution  $q$ , the MMD-based tests are competitive with a simpler derivation in complicated testing  
 650 scenarios [Xu and Matsuda, 2020, 2021], and they can outperform non-kernel based goodness-of-fit  
 651 tests as discussed in Xu and Matsuda [2020].

## 652 C Visual illustrations of the assessment procedures

653 While AgraSSt is illustrated in Figure.1, we provide an additional visualisation emphasising different  
 654 tasks for which AgraSSt can be applied. In Section 4.1 in the main text, we mentioned two features  
 655 of our proposed AgraSSt procedure:



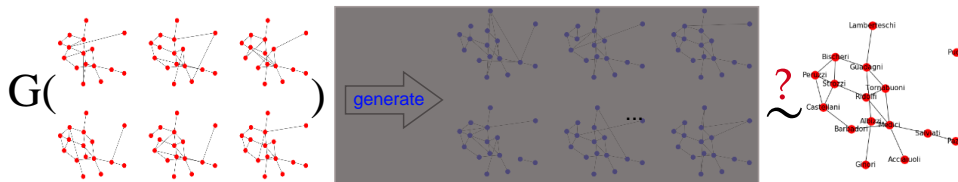


Figure 4: Criticising training quality for generative models.

- 656 1. Regardless of the learning or training procedures (masked in grey), AgraSSt can test a given  
 657 generator  $G$  that is only accessible through its generated samples as shown in Figure.3. In  
 658 this setting, we do not need to know how the generator  $G$  is obtained and the focus is the  
 659 assessment of a particular generator  $G$  itself.
- 660 2. Moreover, we are also interested in understanding the quality and capability of training  
 661 procedures of (deep) generative models. As illustrated in Figure.4, a generator  $G$  is trained  
 662 from the same distribution as the input graph, e.g. ERGMs. The focus in this setting is to  
 663 assess the training procedure of the generative model. (The samples generated are masked in  
 664 grey.) For instance, for  $G$  trained from the Florentine marriage network [Padgett and Ansell,  
 665 1993], we may like to understand whether the generative model can be trained to generate  
 666 graphs that resemble the Florentine marriage network.

667 **D Additional experimental results and discussions**

668 **D.1 Generating reliable samples**

669 To illustrate how AgraSSt can be used to select sample batches, Figure.5 shows three sample batches  
 670 of size 8 for the Karate club network of Zachary [1977], including the corresponding  $p$ -values for  
 671 the displayed sample batches. Here we would expect to detect some community structure in the  
 672 networks; only the sample batch from CELL captures this feature at least to some extent and has  
 673  $p$ -value which would not lead to rejection at the 5% level. This finding chimes with the results from  
 674 Table 2; AgraSSt rejects both GraphRNN and NetGAN as synthetic data generators, but does not  
 675 reject CELL.

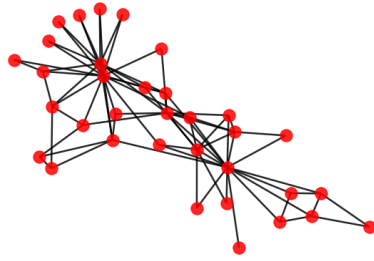
676 **D.2 Additional case study: Padgett’s Florentine network**

677 Padgett’s Florentine network [Padgett and Ansell, 1993]. has 16 vertices and 20 edges; in Xu and  
 678 Reinert [2021] the hypothesis that it is an instance of a  $G(n, p)$  model could not be rejected.

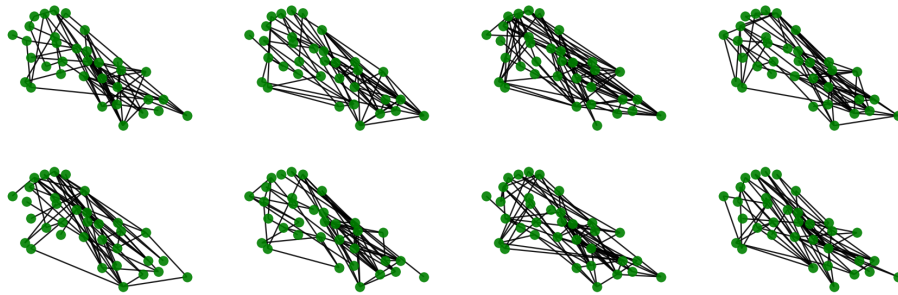
	AgraSSt	Deg	MDdeg	TV_deg
GraphRNN	0.01	0.11	0.26	0.03
NetGAN	0.16	0.18	0.09	0.06
CELL	0.23	0.36	0.69	0.18

Table 3:  $p$ -values for models trained from Florentine marriage network; 100 samples to simulate the null; rejected null at significant level  $\alpha = 0.05$  is marked red.

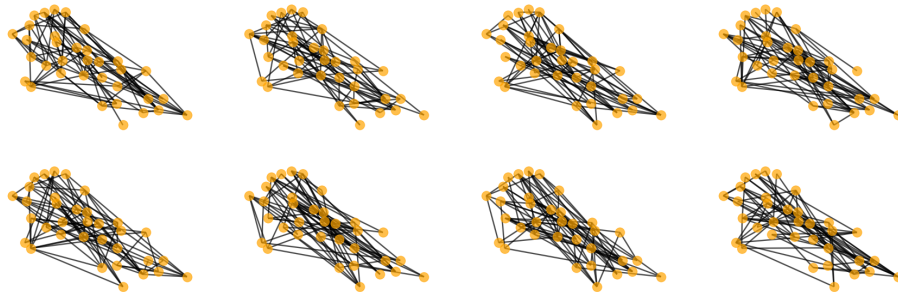
679 The  $p$ -values for different tests are shown in Table.3. The Florentine marriage network has edge  
 680 density  $q = 0.167$ , while the trained CELL has  $\hat{q} = 0.165$  which is a close approximation. GraphRNN  
 681 generates graphs with higher edge density  $\hat{q} = 0.188$ . NetGAN generate samples with  $\hat{q} = 0.176$ ,  
 682 not too different from the null, which is not rejected at  $\alpha = 0.05$ . This is different from what we see  
 683 in the ERGM case above. This discrepancy may arise as the Florentine network is small with  $n = 16$   
 684 and not highly clustered, with average local clustering coefficient 0.191.



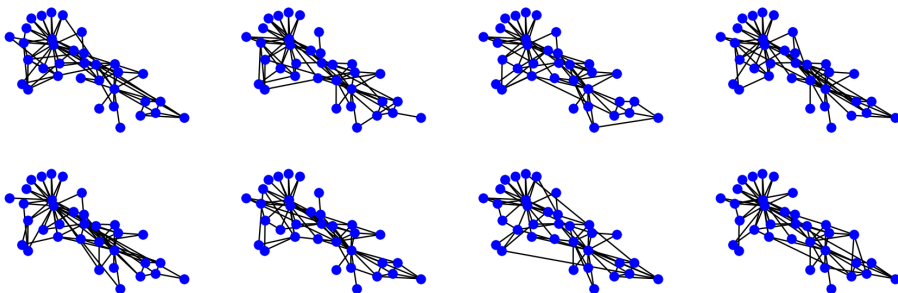
(a) The Karate Club network (vertices in red)



(b) Samples generated from GraphRNN model trained on Karate Club network (vertices in green)



(c) Samples generated from NetGAN model trained on Karate Club network (vertices in orange)



(d) Samples generated from CELL model trained on Karate Club network (vertices in blue) 0.26

Figure 5: The Karate Club network Zachary [1977] and three sample batches of size 8 from different graph generators. The  $p$ -value for GraphRNN samples in (b) is 0.00, for NetGAN samples in (c) the  $p$ -value is 0.01; for CELL samples in (d) the  $p$ -value is 0.26.

685 **Sample batch selection** With CELL being deemed a good generator for the Florentine marriage  
 686 network, we generate a sample batch of size 30 and check the sample quality. Most sample batches  
 687 produce a  $p$ -value above  $\alpha = 0.05$  until the 8th batch, which has  $p$ -value  $0.03 < \alpha$ . AgraSSt would  
 688 recommend not taking this batch. A visual illustration is shown in Figure.6.

689 To investigate these batches, we note that the Florentine marriage network has 3 triangles, while the  
 690 batch being rejected has a significantly lower average number of triangles, namely 1.2. Despite a  
 691 well estimated edge density, this batch produces a low  $p$ -value. This batch, identified by AgraSSt as  
 692 less reliable, may not be very suitable for downstream tasks and it may be better to generate another  
 693 batch instead. In contrast, the batch with  $p$ -value 0.75 has 2.28 triangles on average while the batch  
 694 with  $p$ -value 0.37 has 2.04 triangles on average; these averages are closer to the observed number of  
 695 triangles.

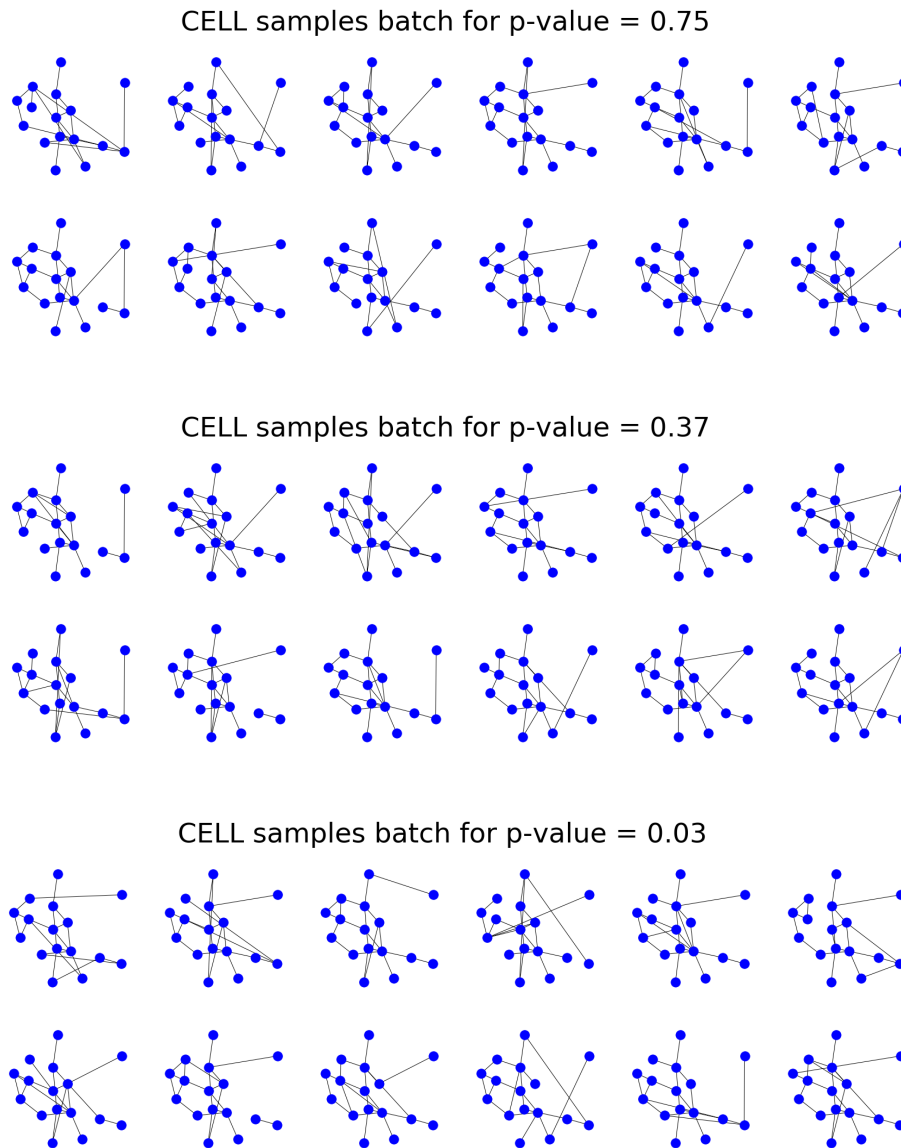


Figure 6: Samples from small size batches generated from CELL trained on the Florentine marriage network, with AgraSSt  $p$ -values. The first two batches would be deemed suitable by AgraSSt, while AgraSSt would not accept the third sample at the 5% significance level.

696 **D.3 Experiments with other network statistics**

697 AgraSSt can incorporate any user-defined network statistics. Table 4 and Table 5 show additional  
 698 results in the settings of Figure 2(b) and Table 1, respectively. As AgraSSt network statistics  $t(x_{-s})$ ,  
 699 we introduce D3, which is based on the multivariate statistics  $(\text{edges}((i,j)), \text{deg}(i), \text{deg}(j))$ , and we  
 700 introduce Tri, which is based on the number of common neighbours of  $i$  and  $j$ . The edge based  
 701 AgraSSt from the main text is added in grey for comparison.

perturbed $\beta_2$	-0.60	-0.40	-0.20	0.00	0.20
AgraSSt_D3	0.93	0.87	0.60	0.06	1.00
AgraSSt_Tri	0.82	0.71	0.35	0.07	1.00
AgraSSt (main)	0.95	0.89	0.68	0.04	1.00

Table 4: Rejection Rate for the setting in Figure 2(b).

Models	GraphRNN	NetGAN	CELL	MC
AgraSSt_D3	0.31	0.66	0.10	0.03
AgraSSt_Tri	0.28	0.32	0.12	0.06
AgraSSt (main)	0.42	0.81	0.05	0.04

Table 5: Rejection Rate for the setting in Table 1.

702 In the Florentine network example, D3 has  $p$ -values 0.04 for GraphRNN, 0.11 for NetGAN, and 0.74  
 703 for CELL. Tri has  $p$ -values 0.02 for GraphRNN, 0.01 for NetGAN, and 0.12 for CELL. Overall, the  
 704 results are mainly comparable to using AgraSSt based on the number of edges, although Tri rejects  
 705 NetGAN for the Florentine marriage network, thus picking up on NetGAN struggling to reproduce  
 706 local clustering.

707 **D.4 Additional discussions on distance-based test statistics**

708 A classical approach for goodness-of-fit testing in ERGMs is the graphical test by Hunter et al. [2008].  
 709 The idea is to simulate sample graphs under the null distribution statistics and create box plots of  
 710 some relevant network statistics; add to these plots the network statistics in the observed network, as  
 711 a solid line for comparison, which is illustrated in Figure.7. The box plot is used to check whether the  
 712 observed network is “very different” from the simulated null samples. This graphical test procedure  
 713 can be translated into Monte Carlo tests. It is natural to adapt such procedure to implicit models from  
 714 which samples can be obtained. Figure.7 plots standard network statistics from Hunter et al. [2008]  
 715 for samples from a fitted  $G(n, p)$  generator (ER Approximate) and a learned GraphRNN generator  
 716 of the Florentine marriage network described in more detail in Appendix D.2. The bold black line  
 717 indicates the distribution of statistics for the Florentine marriage network.

718 The distribution of network statistics is then quantified via Total Variation (TV) distance [Xu and  
 719 Reinert, 2021], based on which goodness-of-fit testing with  $p$ -values can be conducted. We find  
 720 that while the fitted ER generator shows a reasonable fit for all summary statistics, the GraphRNN  
 721 generator does not match the Florentine marriage network very well for dyad-wise shared partners  
 722 and the triad census.

723 **D.5 Efficiency results**

724 Table.6 presents the computational runtime (RT) and the test construction time (CT) for AgraSSt  
 725 and its comparison methods from Section 5.1.1 with the simulation setup as in Section 5.1.2. As a  
 726 measure of accuracy, the variance (Var) of the simulated (or estimated) test statistics under the null  
 727 distribution is also included.

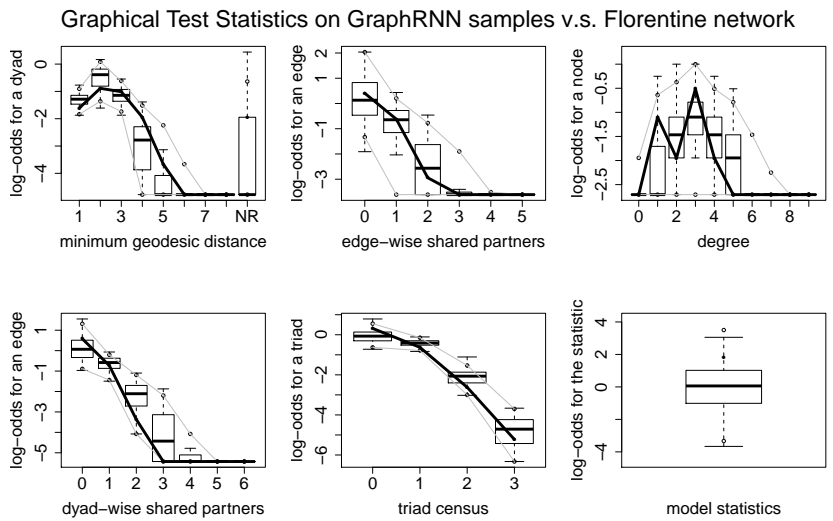
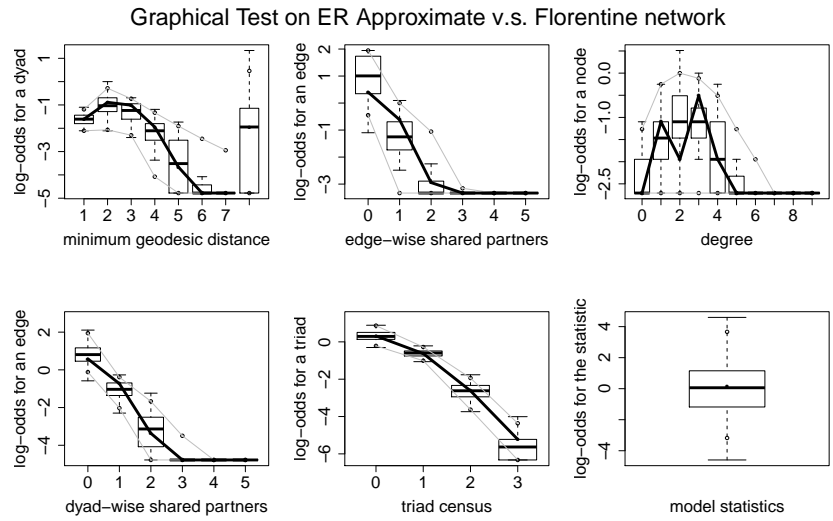


Figure 7: Graphical test illustrations on samples from generators learned on the Florentine marriage network

728 The parameter estimation in **Param** depends on a computationally efficient method which is based  
729 on MPLE [Schmid and Desmarais, 2017] in Eq.(14). **AgraSSt** takes longer to compute mainly due to  
730 the computation of graph kernels, e.g. Weisfeiler-Lehman kernel [Shervashidze et al., 2011]. We note  
731 that for implicit models, the estimation step in AgraSSt relies on generating samples from the model  
732 so that the the computational advantage<sup>4</sup> of the Stein based test over graphical goodness-of-fit tests<sup>5</sup>  
733 reduces compared to gKSS. **MDdeg** is computationally expensive due to the estimation of an inverse  
734 covariance matrix. While providing fast computation and estimation, **Deg** and **Param** sacrifice test  
735 power through a large variance of the test statistics. Estimating the full degree distribution, the total  
736 variation distance method **TV\_deg**, based only on degrees, is competitive with AgraSSt; we recall  
737 that in our simulation results from Section 5.1.2 **TV\_deg** was less powerful than AgraSSt. Here  
738 **MDdeg** is outperformed by the other test statistics.

	AgraSSt	Deg	Param	MDdeg	TV_deg
RT(s)	0.141	0.0006	0.014	0.831	0.002
CT(s)	28.656	0.277	2.963	162.912	0.555
Var	0.23	8.38	1.43	15.84	0.28

Table 6: Computational efficiencies and uncertainty in estimates. RT: runtime for one test; CT: construction time for the test class, including generating 500 samples for relevant estimation and 200 samples for simulating from the null distribution; Var: the estimated variance under the simulated null distribution. Both RT and CT are in seconds.

## 739 D.6 Additional implementation details

740 For GraphRNN, we use batch size 128, epoch 1000 for training, 100 for testing, and learning rate  
741 0.003. For CELL, we use learning rate 0.01, and weight decay 1e-7. For NetGAN, we use batch size  
742 128, epoch 50, generator size and discriminator size both 128, and learning rate 0.0003.

743 We note that training NetGAN [Bojchevski et al., 2018] with the Florentine and with the Karate Club  
744 network may encounter some generator instability and hence early stopping can be useful. Without  
745 early stopping, the training loss for the generator increases during training, although it should be  
746 decreasing. Figure.8 shows the training loss on the generator in NetGAN as well as on the critic  
747 (or discriminator) in NetGAN. Figure.8 plots the loss every 200 training epochs. We can see from  
748 Figure.8(a) that the generator loss starts to be unstable and then increases after 50 points, i.e. 10,000  
749 epochs. Hence we use only 10,000 epochs for training.

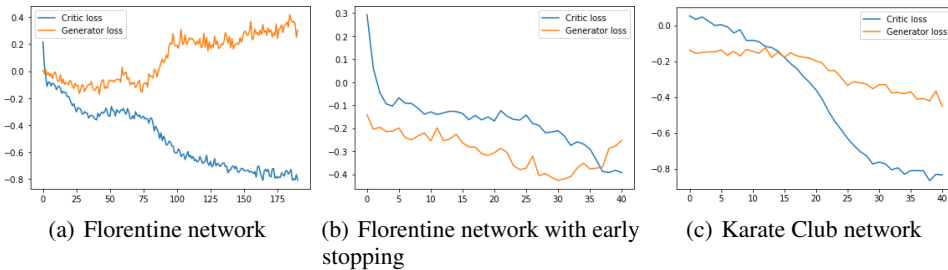


Figure 8: Training loss (y-axis) for NetGAN [Bojchevski et al., 2018]; plotted against every 200 training epochs (x-axis).

<sup>4</sup>These results on gKSS are shown in Supplementary Material D in Xu and Reinert [2021].

<sup>5</sup>The graphical test [Hunter et al., 2008] is computed based on generating a large amount of samples from the null distribution.