# Finding induced subgraphs in scale-free inhomogeneous random graphs

Ellen Cardinaels[1], Johan S.H. van Leeuwaarden[2], and Clara Stegehuis[3]

[1]Eindhoven University of Technology
[2]Tilburg University
[3]Twente University

August 30, 2019

### Abstract

We study the problem of finding a copy of a specific induced subgraph on inhomogeneous random graphs with infinite variance power-law degrees. We provide a fast algorithm that finds a copy of any connected graph $H$ on a fixed number of $k$ vertices as an induced subgraph in a random graph with $n$ vertices. By exploiting the scale-free graph structure, the algorithm runs in $O(nk)$ time for small values of $k$. As a corollary, this shows that the induced subgraph isomorphism problem can be solved in time $O(nk)$ for the inhomogeneous random graph. We test our algorithm on several real-world data sets.

## 1 Introduction

The induced subgraph isomorphism problem asks whether a large graph $G$ contains a connected graph $H$ as an induced subgraph. When $k$ is allowed to grow with the graph size $n$, this problem is NP-hard in general. For example, $k$-clique and $k$ induced cycle, special cases of $H$, are known to be NP-hard [13,19]. For fixed $k$, this problem can be solved in polynomial time $O(n^k)$ by searching for $H$ on all possible combinations of $k$ vertices. Several randomized and non-randomized algorithms exist to improve upon this trivial way of finding $H$ [14,24,26,28].

A second problem we investigate is how to find a subgraph, when we know it exists, more efficiently than the trivial O(nk) algorithm.

On real-world networks, many algorithms were observed to run much faster than predicted by the worst-case running time of algorithms. This may be ascribed to some of the properties that many real-world networks share [4], such as the power-law degree distribution found in many networks [1, 8, 18, 27]. One way of exploiting these power-law degree distributions is to design algorithms that work well on random graphs with power-law degree distributions. For example, finding the largest clique in a network is NP-complete for general networks [19]. However, in random graph models such as the Erdős-Rényi random graph and the inhomogeneous random graph, their specific structures can be exploited to design fixed parameter tractable (FPT) algorithms that efficiently find a clique of size $k$ [10,12] or the largest independent set [15].

In this paper, we study algorithms that are designed to efficiently find subgraphs in the inhomogeneous random graph, a random graph model that can generate graphs with a power-law degree distribution [2, 3, 5, 6, 23, 25]. The inhomogeneous random graph has a densely connected core containing many cliques, consisting of vertices with degrees $\sqrt{n \log(n)}$ and larger. In this densely connected core, the probability of an edge being present is close to one, so that it contains many complete graphs [17]. This observation was exploited in [11] to efficiently determine whether a clique of size $k$ occurs as a subgraph in an inhomogeneous random graph. When searching for *induced* subgraphs however, some edges are required not to be present. Therefore, searching for

induced subgraphs in the entire core is not efficient. We show that a connected subgraph $H$ can be found as an induced subgraph by scanning only vertices that are on the boundary of the core: vertices with degrees proportional to $\sqrt{n}$.

We present an algorithm that first selects the set of vertices with degrees proportional to $\sqrt{n}$, and then randomly searches for $H$ as an induced subgraph on a subset of $k$ of those vertices. The first algorithm we present does not depend on the specific structure of $H$. For general sparse graphs, the best known algorithms to solve subgraph isomorphism on 3 or 4 vertices run in $O(n^{1.41})$ or $O(n^{1.51})$ time with high probability [28]. For small values of $k$, our algorithm finds the desired subgraph on $k$ nodes in linear time with high probability on inhomogeneous random graphs. However, the graph size needs to be very large for our algorithm to perform well. We therefore present a second algorithm that again selects the vertices with degrees proportional to $\sqrt{n}$, and then searches for induced subgraph $H$ in a more efficient way. This algorithm has the same performance guarantee as our first algorithm, but performs much better in simulations.

We test our algorithm on large inhomogeneous random graphs, where it indeed efficiently finds induced subgraphs. We also test our algorithm on real-world network data with power-law degrees. There our algorithm does not perform well, probably due to the fact that the densely connected core of some real-world networks may not be the vertices of degrees at least proportional to $\sqrt{n}$. We then show that a slight modification of our algorithm that looks for induced subgraphs on vertices of degrees proportional to $n^\gamma$ for some other value of $\gamma$ performs better on real-world networks, where the value of $\gamma$ depends on the specific network.

**Notation.** We say that a sequence of events $(\mathcal{E}_n)_{n \geq 1}$ happens with high probability (w.h.p.) if $\lim_{n \to \infty} \mathbb{P}(\mathcal{E}_n) = 1$. Furthermore, we write $f(n) = o(g(n))$ if $\lim_{n \to \infty} f(n)/g(n) = 0$, and $f(n) = O(g(n))$ if $|f(n)|/g(n)$ is uniformly bounded, where $(g(n))_{n \geq 1}$ is nonnegative. Similarly, if $\limsup_{n \to \infty} |f(n)|/g(n) > 0$, we say that $f(n) = \Omega(g(n))$ for nonnegative $(g(n))_{n \geq 1}$. We write $f(n) = \Theta(g(n))$ if $f(n) = O(g(n))$ as well as $f(n) = \Omega(g(n))$.

## 1.1 Model

As a random graph null model, we use the inhomogeneous random graph or hidden variable model $[2, 3, 5, 6, 23, 25]$. Every vertex is equipped with a weight $w$. We assume that the weights are i.i.d. samples from the power-law distribution

$$\mathbb{P}(w > k) = Ck^{1-\tau} \tag{1.1}$$

for some constant $C$ and for $\tau \in (2, 3)$. Two vertices with weights $w$ and $w'$ are connected with probability

$$p(w, w') = \min\left(\frac{ww'}{\mu n}, 1\right), \tag{1.2}$$

where $\mu$ denotes the mean value of the power-law distribution (1.1). Choosing the connection probability in this way ensures that the expected degree of a vertex with weight $w$ is $w$. We denote the degree of vertex $i$ in the inhomogeneous random graph by $D_i$ and its weight by $w_i$.

## 1.2 Algorithms

We now describe two randomized algorithms that determine whether a connected graph $H$ is an induced subgraph in an inhomogeneous random graph and finds the location of such a subgraph if it exists. Algorithm 1 selects the vertices in the inhomogeneous random graph that are on the boundary of the core of the graph: vertices with degrees slightly below $\sqrt{\mu n}$. Then, the algorithm randomly divides these vertices into sets of $k$ vertices. If one of these sets contains $H$ as an induced subgraph, the algorithm terminates and returns the location of $H$. If this is not the case, then the algorithm fails. In the next section, we show that for $k$ small enough, the probability that the algorithm fails is small. This means that $H$ is present as an induced subgraph on vertices that are on the boundary of the core with high probability.

2

Algorithm 1 is similar to the algorithm in [12] designed to find cliques in random graphs. The major difference is that the algorithm to find cliques looks for cliques on all vertices with degrees larger than $\sqrt{f_1\mu n}$ for some function $f_1$. This algorithm is not efficient for detecting subgraphs other than cliques, since vertices with high degrees will be connected with probability close to one.

---

**Algorithm 1:** Finding induced subgraph $H$ (random search)

---

  **Input** : $H = (V_H, E_H)$, $G = (V_G, E_G)$, $\mu$, $f_1 = f_1(n)$, $f_2 = f_2(n)$.
  **Output:** Location of $H$ in $G$ or fail.
**1** Define $n = |V|$, $I_n = [\sqrt{f_1\mu n}, \sqrt{f_2\mu n}]$, set $k = |V_H|$ and $V' = \emptyset$.
**2** **for** $i \in V$ **do**
**3** $\quad$ **if** $D_i \in I_n$ **then** $V' = V' \cup i$;
**4** **end**
**5** Divide the vertices in $V'$ randomly into $\lfloor |V'|/k \rfloor$ sets $S_1, \ldots, S_{\lfloor |V'|/k \rfloor}$.
**6** **for** $j = 1, \ldots, \lfloor |V'|/k \rfloor$ **do**
**7** $\quad$ **if** $H$ *is an induced subgraph on* $S_j$ **then return** location of $H$;
**8** **end**

---

The following theorem gives a bound for the performance of Algorithm 1 for small values of $k$.

**Theorem 1.** *Choose $f_1 = f_1(n) \geq 1/\log(n)$ and $f_2(n)$ such that for some $a < 1, b < 1$, $f_1 < af_2$ and $f_2 < b < 1$ for all $n$. Let $k < \log^{1/3}(n)$. Then, with high probability, Algorithm 1 detects induced subgraph $H$ on $k$ vertices in an inhomogeneous random graph with $n$ vertices and weights distributed as in* (1.1) *in time $O(nk)$.*

Thus, for constant values of $k$, Algorithm 1 finds an instance of $H$ in linear time.

A problem with parameter $k$ is called fixed parameter tractable (FPT) if it can be solved in $f(k)n^{O(1)}$ time for some function $f(k)$, and it is called typical FPT (typFPT) if it can be solved in $f(k)n^{O(1)}$ with high probability [9]. As a corollary of Theorem 1 we obtain that the induced subgraph problem on the inhomogeneous random graph is in typFPT for any subgraph $H$, similarly to the $k$-clique problem on inhomogeneous random graphs [12].

**Corollary 2.** *The induced subgraph problem on the inhomogeneous random graph is in typFPT.*

In theory Algorithm 1 detects any motif on $k$ vertices in linear time for small $k$. However, this only holds for large values of $n$, which can be understood as follows. In Lemma 6, we show that $|V'| = \Theta(n^{(3-\tau)/2})$, thus tending to infinity as $n$ grows large. However, when $n = 10^7$ and $\tau = 2.5$, this means that the size of the set $V'$ is only proportional to $10^{1.75} = 56$ vertices. Therefore, the number of sets $S_j$ constructed in Algorithm 1 is also small. Even though the probability of finding motif $H$ in any such set is proportional to a constant, this constant may be small, so that for finite $n$ the algorithm almost always fails. Thus, for Algorithm 1 to work, $n$ needs to be large enough so that $n^{(3-\tau)/2}$ is large as well.

The algorithm can be significantly improved by changing the search for $H$ on vertices in set $V'$. In Algorithm 2 we propose a search for motif $H$ similar to the Kashtan motif sampling algorithm [20]. Rather than sampling $k$ vertices randomly, it samples one vertex randomly, and then randomly increases the set $S$ by adding vertices in its neighborhood. This already guarantees the vertices in list $S_j$ to be connected, making it more likely for them to form a specific connected motif together. In particular, we expand the list $S_j$ in such a way that the vertices in $S_j$ are guaranteed to form a spanning tree of $H$ as a subgraph. This is ensured by choosing the list $T^H$ that specifies at which vertex in $S_j$ we expand $S_j$ by adding a new vertex. For example, if $k = 4$ and we set $T^H = [1, 2, 3]$ we first add an edge to the first vertex, then we look for a random neighbor of the previously added vertex, and then we add a random neighbor of the third added vertex. Thus, setting $T^H = [1, 2, 3]$ ensures that the set $S_j$ contains a path of length three, whereas setting $T^H = [1, 1, 1]$ ensures that the set $S_j$ contains a star-shaped subgraph. Depending on which subgraph $H$ we are looking for, we can define $T^H$ in such a way that we ensure that the set $S_j$ at least contains a spanning tree of motif $H$ in Step 6 of the algorithm.

3

The selection on the degrees ensures that the degrees are sufficiently high so that probability of finding such a connected set on $k$ vertices is high, as well as that the degrees are sufficiently low to ensure that we do not only find complete graphs because of the densely connected core of the inhomogeneous random graph. The probability that Algorithm 2 indeed finds the desired motif $H$ in any check is of constant order of magnitude, similar to Algorithm 1. Therefore, the performance guarantee of both algorithms is similar. However, for several synthetic and real-world data sets, we show in Section 3 that Algorithm 2 performs much better, since for finite $n$, $k$ connected vertices are more likely to form a motif than $k$ randomly chosen vertices.

---

**Algorithm 2:** Finding induced subgraph $H$ (neighborhood search)

    **Input**   : $H$, $G = (V, E)$, $\mu$, $f_1 = f_1(n)$, $f_2 = f_2(n)$, $s$.
    **Output:** Location of $H$ in $G$ or fail.
**1** Define $n = |V|$, $I_n = [\sqrt{f_1 \mu n}, \sqrt{f_2 \mu n}]$ and set $V' = \emptyset$.
**2** **for** $i \in V$ **do**
**3**    |   **if** $D_i \in I_n$ **then**  $V' = V' \cup i$;
**4** **end**
**5** Let $G'$ be the induced subgraph of $G$ on vertices $V'$.
**6** Set $T^H$ consistently with motif $H$ .
**7** **for** $j=1,\ldots,s$ **do**
**8**    |   Pick a random vertex $v \in V'$ and set $S_j = \{v\}$.
**9**    |   **while** $|S_j| \neq k$ **do**
**10**    |   |   Pick a random $v' \in N_{G'}(S_j[T^H[j]]) : v' \notin S_j$
**11**    |   |   Add $v'$ to $S_j$.
**12**    |   **end**
**13**    |   **if** *H is an induced subgraph on $S_j$* **then return** location of $H$;
**14** **end**

---

The following theorem shows that indeed Algorithm 2 has similar performance guarantees as Algorithm 1.

**Theorem 3.** *Choose $f_1 = f_1(n) \geq 1/\log(n)$ and $f_1 < f_2 < 1$. Choose $s = \Omega(n^\alpha)$ for some $0 < \alpha < 1$, such that $s \leq n/k$. Then, Algorithm 2 detects induced subgraph $H$ on $k < \log^{1/3}(n)$ vertices on an inhomogeneous random graph with $n$ vertices and weights distributed as in (1.1) in time $O(nk)$ with high probability.*

The proofs of Theorem 1 and 3 rely on the fact that for small $k$, any subgraph on $k$ vertices is present in $G'$ with high probability. This means that after the degree selection step of Algorithms 1 and 2, for small $k$, any motif finding algorithm can be used to find motif $H$ on the remaining graph $G'$, such as the Grochow-Kellis algorithm [14], the MAvisto algorithm [26] or the MODA algorithm [24]. In the proofs of Theorem 1 and 3, we show that $G'$ has $\Theta(n^{(3-\tau)/2})$ vertices with high probability. Thus, the degree selection step reduces the problem of finding a motif $H$ on $n$ vertices to finding a motif on a graph with $\Theta(n^{(3-\tau)/2})$ vertices, significantly reducing the running time of the algorithms.

## 2   Proof of Theorems 1 and 3

We prove Theorem 1 using two lemmas. The first lemma relates the degrees of the vertices to their weights. The connection probabilities in the inhomogeneous random graph depend on the weights of the vertices. In Algorithm 1, we select vertices based on their degrees instead of their unknown weights. The following lemma shows that the weights of the vertices in $V'$ are close to their degrees.

4

**Lemma 4** (Degrees and weights). *Fix $\varepsilon > 0$, and define $J_n = [(1-\varepsilon)\sqrt{f_1\mu n}, (1+\varepsilon)\sqrt{f_2\mu n}]$. Then, for some $K > 0$,*

$$\mathbb{P}\left(\exists i \in V' : w_i \notin J_n\right) \leq Kn \exp\left(-\varepsilon^2\sqrt{\mu n}\min\left(\frac{\sqrt{f_1}}{1-\varepsilon}, \frac{\sqrt{f_2}}{1+\varepsilon}\right)/2\right). \tag{2.1}$$

*Proof.* Fix a vertex $i \in V$. Then,

$$\begin{aligned}
\mathbb{P}\left(w_i < (1-\varepsilon)\sqrt{f_1\mu n},\ D_i \in I_n\right) &= \frac{\mathbb{P}\left(D_i \in I_n \mid w_i < (1-\varepsilon)\sqrt{f_1\mu n}\right)}{\mathbb{P}\left(w_i < (1-\varepsilon)\sqrt{f_1\mu n}\right)} \\
&\leq \frac{\mathbb{P}\left(D_i > \sqrt{f_1\mu n} \mid w_i = (1-\varepsilon)\sqrt{f_1\mu n}\right)}{1 - C((1-\varepsilon)\sqrt{f_1\mu n})^{1-\tau}} \\
&\leq K_1\mathbb{P}\left(D_i > \sqrt{f_1\mu n} \mid w_i = (1-\varepsilon)\sqrt{f_1\mu n}\right), \tag{2.2}
\end{aligned}$$

for some $K_1 > 0$.

Here the first inequality follows because the probability that a vertex with weight $w_1$ has degree at least $\sqrt{f_1\mu n}$ is larger than the probability that a vertex of weight $w_2$ has degree at least $\sqrt{f_1\mu n}$ when $w_1 > w_2$. Conditionally on the weights, $D_i$ is the sum of $n-1$ independent indicators indicating the presence of an edge between vertex $i$ and the other vertices and that $\mathbb{E}[D_i] = w_i$. Therefore, by the Chernoff bound

$$\mathbb{P}\left(D_i > w_i(1+\delta)\right) \leq \exp\left(-\delta^2 w_i/2\right). \tag{2.3}$$

Therefore, choosing $\delta = \varepsilon/(1-\varepsilon)$ yields

$$\mathbb{P}\left(D_i > \sqrt{f_1\mu n} \mid w_i = (1-\varepsilon)\sqrt{f_1\mu n}\right) \leq \exp\left(-\frac{\varepsilon^2\sqrt{f_1\mu n}}{2(1-\varepsilon)}\right)(1 + o(1)). \tag{2.4}$$

Combining this with (2.2) and taking the union bound over all vertices then results in

$$\mathbb{P}\left(\exists i : D_i \in I_n, w_i < (1-\varepsilon)\sqrt{f_1\mu n}\right) \leq K_2 n \exp\left(-\frac{\varepsilon^2}{2(1-\varepsilon)}\sqrt{f_1\mu n}\right), \tag{2.5}$$

for some $K_2 > 0$. Similarly,

$$\mathbb{P}\left(\exists i : D_i \in I_n, w_i > (1+\varepsilon)\sqrt{f_2\mu n}\right) \leq K_3 n \exp\left(-\frac{\varepsilon^2}{2(1+\varepsilon)}\sqrt{f_2\mu n}\right), \tag{2.6}$$

for some $K_3 > 0$, which proves the lemma. $\qquad\square$

**Lemma 5** (Weights and degrees). *Fix $\varepsilon > 0$, sufficiently small so that $\tilde{J}_n = [(1+\varepsilon)\sqrt{f_1\mu n}, (1-\varepsilon)\sqrt{f_2\mu n}]$ is a non-empty interval. Then, for some $K > 0$,*

$$\mathbb{P}\left(\exists i : w_i \in \tilde{J}_n, i \notin V'\right) \leq Kn \exp\left(-\varepsilon^2\sqrt{\mu n}\min\left(\frac{\sqrt{f_1}}{1-\varepsilon}, \frac{\sqrt{f_2}}{1+\varepsilon}\right)/2\right). \tag{2.7}$$

*Proof.* Fix a vertex $i$ with $w_i \in \tilde{J}_n$. Then,

$$\mathbb{P}\left(D_i < \sqrt{f_1\mu n} \mid w_i \in \tilde{J}_n\right) \leq \mathbb{P}\left(D_i < \sqrt{f_1\mu n} \mid w_i = (1+\varepsilon)\sqrt{f_1\mu n}\right). \tag{2.8}$$

Similarly to (2.4),

$$\mathbb{P}\left(D_i < \sqrt{f_1\mu n} \mid w_i = (1+\varepsilon)\sqrt{f_1\mu n}\right) \leq \exp\left(-\frac{\varepsilon^2\sqrt{f_1 n}}{2(1+\varepsilon)}\right), \tag{2.9}$$

so that

$$\mathbb{P}\left(\exists i : w_i \in \tilde{J}_n, D_i < \sqrt{f_1\mu n}\right) \leq K_1 n \exp\left(-\frac{\varepsilon^2\sqrt{f_1 n}}{2(1+\varepsilon)}\right). \tag{2.10}$$

5

Similarly,

$$\mathbb{P}\left(\exists i : w_i \in \tilde{J}_n, D_i > \sqrt{f_2\mu n}\right) \leq K_2 n \exp\left(-\frac{\varepsilon^2 \sqrt{f_2 n}}{2(1-\varepsilon)}\right). \tag{2.11}$$

□

The second lemma shows that after deleting all vertices with degrees outside of $I_n$ defined in Step 1 of Algorithm 1, still polynomially many vertices remain with high probability.

**Lemma 6.** *Polynomially many nodes remain. There exists $\gamma, \gamma' > 0$ such that*

$$\mathbb{P}\left(|V'| < \gamma n^{(3-\tau)/2}\right) \leq 2\exp\left(-\Theta(n^{(3-\tau)/2})\right) \tag{2.12}$$

*and*

$$\mathbb{P}\left(|V'| > \gamma' n^{(3-\tau)/2} \log^{\tau-1}(n)\right) \leq 2\exp\left(-\Theta(n^{(3-\tau)/2})\right) \tag{2.13}$$

*Proof.* Let $\mathcal{E}$ denote the event that all vertices with $w_i \in \tilde{J}_n$ satisfy $i \in V'$ for some $\varepsilon > 0$, with $\tilde{J}_n$ as in Lemma 5. Let $W'$ be the set of all vertices with weights in $\tilde{J}_n$. Conditioned on the event $\mathcal{E}$, any vertex in $W'$ is also in $V'$ so that $|W'| \leq |V'|$. Then, by Lemma 4

$$\mathbb{P}\left(|V'| < \gamma n^{(3-\tau)/2}\right) \leq \mathbb{P}\left(|W'| < \gamma n^{(3-\tau)/2}\right) + Kn\exp\left(-\varepsilon^2\sqrt{\mu n}\min\left(\frac{\sqrt{f_1}}{1-\varepsilon}, \frac{\sqrt{f_2}}{1+\varepsilon}\right)/2\right) \tag{2.14}$$

Furthermore,

$$\mathbb{P}\left(w_i \in \tilde{J}_n\right) = C((1+\varepsilon)\sqrt{f_1\mu n})^{1-\tau} - C((1-\varepsilon)\sqrt{f_2\mu n})^{1-\tau} \geq c_1(\sqrt{n})^{1-\tau} \tag{2.15}$$

for some constant $c_1 > 0$ when $\varepsilon$ is sufficiently small, because $f_1 < af_2$ for some constant $a < 1$ by assumption. Thus, each of the $n$ vertices is in set $W'$ independently with probability at least $c_1(\sqrt{\mu n})^{1-\tau}$. Choose $0 < \gamma < c_1$. Applying the multiplicative Chernoff bound then shows that

$$\mathbb{P}\left(|W'| < \gamma n^{(3-\tau)/2}\right) \leq \exp\left(-\frac{(c_1-\gamma)^2}{2c_1}n^{(3-\tau)/2}\right), \tag{2.16}$$

which proves the first part of the lemma together with (2.17).

Let $\mathcal{E}'$ denote the event that all vertices $i \in V'$ satisfy $w_i \in J_n$ for some $\varepsilon > 0$, with $J_n$ as in Lemma 4. Let $U'$ be the set of all vertices with weights in $J_n$. On the event $\mathcal{E}'$, any vertex in $V'$ is also in $U'$ so that $|U'| \geq |V'|$. Then, by Lemma 4

$$\mathbb{P}\left(|V'| > \gamma' n^{(3-\tau)/2}\right) \leq \mathbb{P}\left(|U'| > \gamma' n^{(3-\tau)/2}\right) + Kn\exp\left(-\varepsilon^2\sqrt{\mu n}\min\left(\frac{\sqrt{f_1}}{1-\varepsilon}, \frac{\sqrt{f_2}}{1+\varepsilon}\right)/2\right). \tag{2.17}$$

Furthermore, for some $c_2 > 0$,

$$\mathbb{P}\left(w_i \in J_n\right) = C((1-\varepsilon)\sqrt{f_1\mu n})^{1-\tau} - C((1+\varepsilon)\sqrt{f_2\mu n})^{1-\tau} \leq c_2(\sqrt{n}/\log(n))^{1-\tau}, \tag{2.18}$$

where we used that $f_1 \geq 1/\log(n)$. Similarly to (2.16),

$$\mathbb{P}\left(|U'| > \gamma' n^{(3-\tau)/2} \log^{\tau-1}(n)\right) \leq \exp\left(-\frac{(c_2-\gamma')^2}{2c_2}\log^{\tau-1}(n)n^{(3-\tau)/2}\right), \tag{2.19}$$

which proves the second part of lemma.

□

We now use these lemmas to prove Theorem 1.

*Proof of Theorem 1.* We condition on the event that $V'$ is of polynomial size (Lemma 6) and that the weights are within the constructed lower and upper bounds (Lemma 4), since both events occur with high probability. This bounds the edge probability between any pair of nodes $i$ and $j$ in $V'$ as

$$p_{ij} < \min\left(\frac{(1+\varepsilon)\sqrt{f_2 \mu n}(1+\varepsilon)\sqrt{f_2 \mu n}}{\mu n}, 1\right) = f_2(1+\varepsilon)^2. \tag{2.20}$$

Because $f_2 < b < 1$ for some constant $b$, $p_{ij} \leq p_+ < 1$ for some constant $p_+$ if we choose $\varepsilon$ small enough. Similarly,

$$p_{ij} > \min\left(\frac{(1-\varepsilon)^2\sqrt{f_1 \mu n}^2}{\mu n}, 1\right) = \Theta\left(\frac{1}{\log(n)}\right), \tag{2.21}$$

by our choice of $f_1$, so that $p_{ij} \geq c_2/\log(n) =: p_-$ for some constant $c_2$. Let $E := |E_H|$ be the number of edges in $H$. We upper bound the probability of not finding $H$ in one of the partitions of size $k$ of $V'$ as $1 - p_-^E(1-p_+)^{\binom{k}{2}-E}$. Since all partitions are disjoint we can upper bound the probability of not finding $H$ in any of the partitions as

$$\mathbb{P}\left(H \text{ not in any set of partitions}\right) \leq \left(1 - p_-^E(1-p_+)^{\binom{k}{2}-E}\right)^{\left\lfloor\frac{|V'|}{k}\right\rfloor}. \tag{2.22}$$

Using that $E \leq k^2$, $\binom{k}{2} - E \leq k^2$ and that $1 - x \leq e^{-x}$ results in

$$\mathbb{P}\left(H \text{ not in the partitions}\right) \leq \exp\left(-p_-^{k^2}(1-p_+)^{k^2}\left\lfloor\frac{|V'|}{k}\right\rfloor\right). \tag{2.23}$$

Since $|V'| = \Omega\left(n^{\frac{3-\tau}{2}}\right)$, $\lceil|V'|/k\rceil \geq dn^{\frac{3-\tau}{2}}/k$ for some constant $d > 0$. We fill in the expressions for $p_-$ and $p_+$, with $c_3 > 0$ a constant

$$\mathbb{P}\left(H \text{ not in the partitions}\right) \leq \exp\left(-\frac{dn^{\frac{3-\tau}{2}}}{k}\left(\frac{c_3}{\log n}\right)^{k^2}\right). \tag{2.24}$$

Now apply that $k \leq \log^{\frac{1}{3}}(n)$. Then

$$\begin{aligned}\mathbb{P}\left(H \text{ not in the partitions}\right) &\leq \exp\left(-\frac{dn^{\frac{3-\tau}{2}}}{\log^{\frac{1}{3}} n}\left(\frac{c_3}{\log n}\right)^{\log^{\frac{2}{3}} n}\right) \\ &\leq \exp\left(-dn^{\frac{3-\tau}{2}-o(1)}\right).\end{aligned} \tag{2.25}$$

Hence, the inner expression grows polynomially such that the probability of not finding $H$ in one of the partitions is negligibly small. The running time of the partial search is given by

$$\frac{|V'|}{k}\binom{k}{2} \leq \frac{n}{k}\binom{k}{2} \leq nk \leq ne^{k^4}, \tag{2.26}$$

which concludes the proof for $k \leq \log^{1/3}(n)$. $\qquad\square$

*Proof of Corollary 2.* If $k > \log^{\frac{1}{3}}(n)$, $n < e^{k^3}$, so that the time it takes to solve the subgraph isomorphism problem is bounded by a function of $k$.

For $k \leq \log^{\frac{1}{3}}(n)$, Theorem 1 shows that the induced subgraph isomorphism problem can be solved in time $nk \leq ne^{k^4}$. Thus, with high probability the induced subgraph isomorphism problem can be solved in $ne^{k^4}$ time, which proves that it is in typFPT. $\qquad\square$

*Proof of Theorem 3.* The proof of Theorem 3 is very similar to the proof of Theorem 1. The only way Algorithm 2 differs from Algorithm 1 is in the selection of the sets $S_j$. As in the previous theorem, we condition on the event that $\gamma n^{(3-\tau)/2} \leq |V'| \leq \gamma' n^{(3-\tau)/2} \log^{\tau-1}(n)$ (Lemma 6) and that the weights of the vertices in $G'$ are bounded as in Lemma 4.

The graph $G' = (V', E')$ constructed in Step 5 of Algorithm 2 then consists of $\Theta(n^{(3-\tau)/2})$ vertices. Furthermore, by the bound (2.21) on the connection probabilities of all vertices in $G'$, the expected degree of a vertex $i$ in $G'$, $D_{i,G'}$, satisfies $\mathbb{E}[D_{i,G'}] = \Omega(n^{(3-\tau)/2}/\log(n))$. We can use similar arguments as in Lemma 4 to show that $D_{i,G'} = \Omega(n^{(3-\tau)/2}/\log(n))$ with high probability for all vertices in $G'$. Since $G'$ consists of $O(n^{(3-\tau)/2}\log(n)^{\tau-1})$ vertices, $D_{i,G'} = O(n^{(3-\tau)/2}\log(n)^{\tau-1})$ as well. This means that for $k < \log^{\frac{1}{3}}(n)$, Steps 8-11 are able to find a connected subgraph on $k$ vertices with high probability.

We now compute the probability that $S_j$ is disjoint with the previous $j-1$ constructed sets. The number of vertices in sets $S_{j-1}, \ldots, S_1$ is bounded by $(j-1)k$. The probability that the first vertex does not overlap with the previous sets is therefore bounded by $1 - (j-1)k/|V'|$, since that vertex is chosen uniformly at random. The second vertex is chosen in a size-biased manner, since it is chosen by following a random edge. The probability that vertex $i$ is added can therefore be bounded as

$$\mathbb{P}\,(\text{vertex } i \text{ is added}) = \frac{D_{i,G'}}{2|E'|} \leq \frac{M\log^{\tau-1}(n)}{|V'|} \tag{2.27}$$

for some constant $M > 0$ by the conditions on the degrees. Therefore, the probability that $S_j$ does not overlap with one of the previously chosen (at most $(j-1)k$) vertices can be bounded from below by

$$\mathbb{P}\,(S_j \text{ does not overlap with previous sets}) \geq \left(1 - \frac{k(j-1)}{|V'|}\right)\left(1 - \frac{Mk(j-1)\log^{\tau-1}(n)}{|V'|}\right)^{k-1}. \tag{2.28}$$

Thus, the probability that all $j$ sets do not overlap can be bounded as

$$\mathbb{P}\,(S_j \cap S_{j-1} \cdots \cap S_1 = \emptyset) \geq \left(1 - \frac{Mk(j-1)\log^{\tau-1}(n)}{|V'|}\right)^{(j-1)k}, \tag{2.29}$$

which tends to one when $jk = o(n^{(3-\tau)/4})$. Let $s_{\mathrm{dis}}$ denote the number of disjoint sets out of the $s$ sets constructed in Algorithm 2. Then, when $s = \Omega(n^\alpha)$ for some $\alpha > 0$, $s_{\mathrm{dis}} > n^\beta$ for some $\beta > 0$ with high probability, because $k < \log^{1/3}(n)$.

The probability that $H$ is present as an induced subgraph is bounded similarly as in Theorem 1. We already know that $k-1$ edges are present. For all other $E-(k-1)$ edges of $H$, and all $\binom{k}{2} - E$ edges that are not present in $H$, we can again use (2.20) and (2.21) to bound on the probability of edges being present or not being present between vertices in $V'$. Therefore, we can bound the probability that $H$ is not found similarly to (2.23) as

$$\mathbb{P}\,(H \text{ not in the partitions}) \leq \mathbb{P}\,(H \text{ not in the disjoint partitions})$$

$$\leq \exp\left(-p_-^{k^2}(1-p_+)^{k^2}s_{\mathrm{dis}}\right).$$

Because $s_{\mathrm{dis}} > n^\beta$ for some $\beta > 0$, this term tends to zero exponentially. The running time of the partial search can be bounded similarly to (2.26) as

$$s\binom{k}{2} \leq sk^2 = O(nk), \tag{2.30}$$

where we used that $s \leq n/k$. $\qquad\square$

## 3   Experimental results

Figure 1 shows the success rate of Algorithm 1, defined as the fraction of times Algorithm 1 succeeds in finding a cycle of size $k$ in an inhomogeneous random graph on $10^7$ vertices. Even
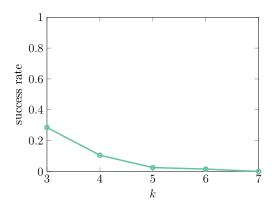
Figure 1: The fraction of times step 7 in Algorithm 1 succeeds to find a cycle of length $k$ in an inhomogeneous random graph with $N = 10^7$, $\tau = 2.5$, averaged over 500 network samples with $f_1 = 1/\log(n)$ and $f_2 = 0.9$.

though for large $n$ Algorithm 1 should find an instance of a cycle of size $k$ in step 7 of the algorithm with high probability, we see that Algorithm 1 never succeeds in finding one of size 7. The success rate of Algorithm 1 on smaller cycles is also far away from 1, because of the finite size effects discussed before.

Figure 2a also plots the fraction of times Algorithm 2 succeeds to find a cycle. We set the parameter $s = 10000$ so that the algorithm fails if the algorithm does not succeed to detect motif $H$ after executing step 13 of Algorithm 2 10000 times. Because $s$ gives the number of attempts to find $H$, increasing $s$ may increase the success probability of Algorithm 2 at the cost of a higher running time. However, in Figure 2b, for small values of $k$, the mean number of times Step 13 is executed when the algorithm succeeds is much lower than 10000, so that increasing $s$ in this experiment probably only has a small effect on the success probability.

Algorithm 2 with $f_1 = 1/\log(n)$ and $f_2 = 0.9$ in line with Theorem 3 outperforms Algorithm 1. Figure 2b also shows that the number of attempts needed to detect a cycle of length $k$ is small for $k \leq 6$. For larger values of $k$ the number of attempts increases. This can again be ascribed to the finite size effects that cause the set $V'$ to be small, so that large motifs may not be present on vertices in set $V'$.

We also plot the success probability when using different values of the functions $f_1$ and $f_2$, outside the window where Theorem 3 holds. When $f_2 = \infty$, all vertices of degree at least $\sqrt{f_1 \mu n}$ are included in $V'$, as in [11]. In this setting, the success probability of the algorithm decreases. This is because the set $V'$ now contains many high degree vertices that are much more likely to form clique motifs than cycles or other connected motifs on $k$ vertices. This makes $f_1 = 1/\log(n), f_2 = \infty$ a very efficient setting for detecting clique motifs [11]. For the cycle motif however, Figure 2b shows that more checks are needed before a cycle is detected, and in some cases the cycle is not detected at all.

The setting $f_1 = 0$ and $f_2 = \infty$ is equivalent to including all vertices in $|V'|$. This is also less efficient than the setting of Theorem 3, as Figure 2a shows. In this situation, the number of attempts needed to find a cycle of length $k$ is larger than for Algorithm 2 for $k \leq 6$.

Most results in this section are about finding cycles, since cycles easily scale in $k$ (larger cycles). However, Algorithm 2 with parameters as in Theorem 3 also helps to find other subgraphs than cycles. Table 1 presents the probability that the algorithm detects the subgraph of Figure 3. Indeed, the probability that Algorithm 2 with parameters chosen as in Theorem 3 succeeds is much larger than the probability that the algorithm succeeds on the whole data set.
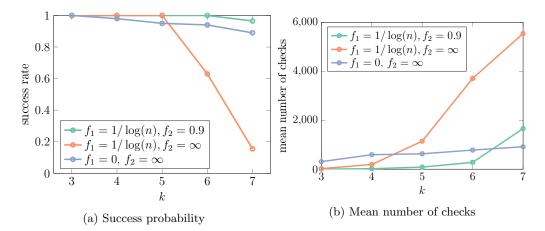
(a) Success probability

(b) Mean number of checks

Figure 2: Results of Algorithm 2 on an inhomogeneous random graph with $N = 10^7$, $\tau = 2.5$ for detecting cycles of length $k$, using $s = 10000$. The values are averaged over 500 generated networks.

| | $f_1 = 1/\log(n), f_2 = 0.9$ | $f_1 = 1/\log(n), f_2 = \infty$ | $f_1 = 0, f_2 = \infty$ |
|---|---|---|---|
| **Success rate** | 0.47 | 0 | 0 |

Table 1: Success rate of Algorithm 2 on the subgraph of Figure 3 for $N = 10^6$, $\tau = 2.5$, $s = 10000$ over 500 generated networks.

## 3.1 Real network data

We now check Algorithm 2 on four real-world networks with power-law degrees: a Wikipedia communication network [21], the Gowalla social network [21], the Baidu online encyclopedia [22] and the Internet on the autonomous systems level [21]. Table 2 presents several statistics of these scale-free data sets. Figure 4 shows the fraction of runs where Algorithm 2 finds a cycle as an induced subgraph. We see that for the Wikipedia social network in Figure 4a, Algorithm 2 with $f_1 = 1/\log(n)$ and $f_2 = 0.9$ in line with Theorem 3 is more efficient than looking for cycles among all vertices in the network ($f_1 = 0, f_2 = \infty$). For the Baidu online encyclopedia in Figure 4c however, we see that Algorithm 2 with $f_1 = 1/\log(n)$ and $f_2 = 0.9$ performs much worse than looking for cycles among all possible vertices. In the other two network data sets in Figures 4b and 4d the performance on the reduced vertex set and the original vertex set is almost the same. Figure 5 shows that in general, Algorithm 2 with settings as in Theorem 3 indeed seems to finish in fewer steps than when using the full vertex set. However, as Figure 5c shows, for larger values of $k$ the algorithm fails almost always.

These results show that while Algorithm 2 with $f_1, f_2$ in line with Theorem 3 is efficient on inhomogeneous random graphs, it may not always be efficient on real-world data sets. This is not surprising, because there is no reason why in real-world data the vertices of degrees proportional to $\sqrt{n}$ should behave like an Erdős-Rényi random graph, like in the inhomogeneous random graph. Thus, in terms of subgraphs, the inhomogeneous random graph and real-world network data differ



Figure 3: The subgraph corresponding to the results in Table 1.

|  | $n$ | $E$ | $\tau$ |
|---|---|---|---|
| **Wikipedia** | 2,394,385 | 5,021,410 | 2.46 |
| **Gowalla** | 196,591 | 950,327 | 2.65 |
| **Baidu** | 2,141,300 | 17,794,839 | 2.29 |
| **AS-Skitter** | 1,696,415 | 11,095,298 | 2.35 |

Table 2: Statistics of the data sets: the number of vertices $n$, the number of edges $E$, and the power-law exponent $\tau$ fitted by the method of [7].

significantly.

We therefore investigate whether selecting vertices with degrees in $I_n = [(\mu n)^\gamma / \log(n), (\mu n)^\gamma]$ for some other value of $\gamma$ in Algorithm 2 leads to a better performance. Figure 4 and 5 show for every data set one particular value of $\gamma$ that works well. For the Gowalla, Wikipedia and Autonomous systems network, this leads to a faster algorithm to detect cycles. In these examples, the success probability of the algorithm is similar to the success probability on the full data set, but Figure 5 shows that it finds the cycle much faster than the algorithm on the full data set. Only for the Baidu network other values of $\gamma$ do not improve upon randomly selecting from all vertices. This indicates that for most networks, cycles do appear mostly on degrees with specific orders of magnitude, making it possible to sample these cycles faster. Unfortunately, these orders of magnitude may be different for different networks. Across all four networks, the best value of $\gamma$ seems to be smaller than the value of 0.5 that is optimal for the inhomogeneous random graph.

In these experiments, we tested the values $\gamma = 0.1, 0.2, 0.3, 0.4, 0.5$, and we present in Figure 4 the values of $\gamma$ that worked best for each data set. However, it would be useful to be able to select the best value of $\gamma$ without trying several values at first. For example, it may be possible to relate $\gamma$ to the degree-exponent $\tau$, or to a specific quantile of the degree sequence. Finding efficient methods to estimate $\gamma$ from a data set is an interesting question for further work.

## 4  Conclusion

We presented an algorithm which solves the induced subgraph problem on inhomogeneous random graphs with infinite variance power-law degrees in time $O(ne^{k^4})$ with high probability as $n$ grows large. This algorithm is based on the observation that for fixed $k$, any subgraph is present on $k$ vertices with degrees slightly smaller than $\sqrt{\mu n}$ with positive probability. Therefore, the algorithm first selects vertices with those degrees, and then uses a random search method to look for the induced subgraph on those vertices.

We show that this algorithm performs well on simulations of inhomogeneous random graphs. Its performance on real-world data sets varies for different data sets. This indicates that the degrees that contain the most induced subgraphs of size $k$ in real-world networks may not be close to $\sqrt{n}$. We then show that on these data sets, it may be more efficient to find induced subgraphs on degrees proportional to $n^\gamma$ for some other value of $\gamma$. The value of $\gamma$ may be different for different networks.

Our algorithm exploits that induced subgraphs are likely formed among $\sqrt{\mu n}$-degree vertices. However, certain subgraphs may occur more frequently on vertices of other degrees [16]. For example, star-shaped subgraphs on $k$ vertices appear more often on one vertex with degree much higher than $\sqrt{\mu n}$ corresponding to the middle vertex of the star, and $k - 1$ lower-degree vertices corresponding to the leafs of the star [16]. An interesting open question is whether there exist better degree-selection steps for specific subgraphs than the one used in Algorithms 1 and 2.
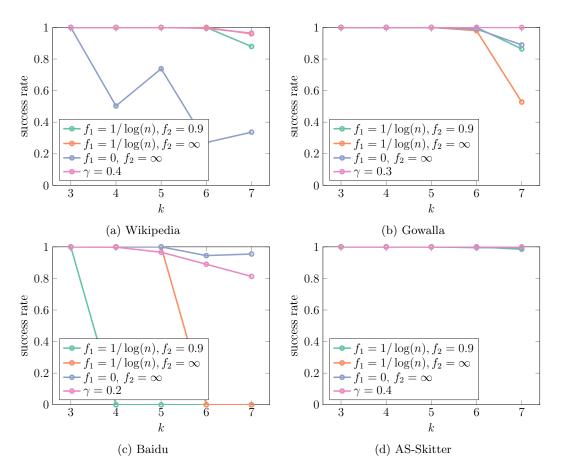
(a) Wikipedia

(b) Gowalla

(c) Baidu

(d) AS-Skitter

Figure 4: The fraction of times Algorithm 2 succeeds to find a cycle on four large network data sets for detecting cycles of length $k$, using $s = 10000$. The pink line uses Algorithm 2 on vertices of degrees in $I_n = [(\mu n)^\gamma / \log(n), (\mu n)^\gamma]$. The values are averaged over 500 runs of Algorithm 2.

# References

[1] Albert, R., Jeong, H., Barabási, A.L.: Internet: Diameter of the world-wide web. Nature **401**(6749) (1999) 130–131

[2] Boguñá, M., Pastor-Satorras, R.: Class of correlated random networks with hidden variables. Phys. Rev. E **68** (2003) 036112

[3] Bollobás, B., Janson, S., Riordan, O.: The phase transition in inhomogeneous random graphs. Random Structures & Algorithms **31**(1) (2007) 3–122

[4] Brach, P., Cygan, M., Łacki, J., Sankowski, P.: Algorithmic complexity of power law networks. In: Proceedings of the Twenty-seventh Annual ACM-SIAM Symposium on Discrete Algorithms. SODA '16, Philadelphia, PA, USA, Society for Industrial and Applied Mathematics (2016) 1306–1325

[5] Britton, T., Deijfen, M., Martin-Löf, A.: Generating simple random graphs with prescribed degree distribution. J. Stat. Phys. **124**(6) (2006) 1377–1397
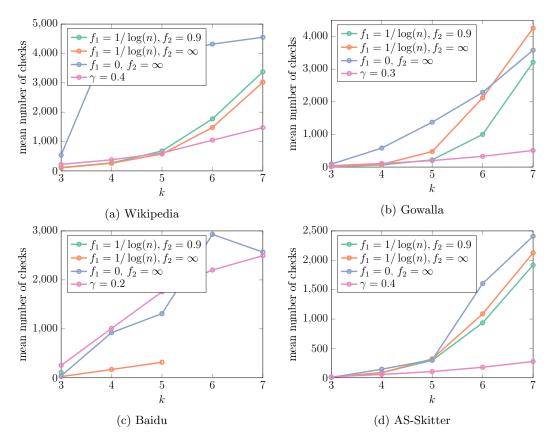
Figure 5: The number of times step 12 of Algorithm 2 is invoked when the algorithm does not fail on four large network data sets for detecting cycles of length $k$, using $s = 10000$. The pink line uses Algorithm 2 on vertices of degrees in $I_n = [(\mu n)^\gamma / \log(n), (\mu n)^\gamma]$. The values are averaged over 500 runs of Algorithm 2.

[6] Chung, F., Lu, L.: The average distances in random graphs with given expected degrees. Proc. Natl. Acad. Sci. USA **99**(25) (2002) 15879–15882 (electronic)

[7] Clauset, A., Shalizi, C.R., Newman, M.E.J.: Power-law distributions in empirical data. SIAM Rev. **51**(4) (2009) 661–703

[8] Faloutsos, M., Faloutsos, P., Faloutsos, C.: On power-law relationships of the internet topology. In: ACM SIGCOMM Computer Communication Review. Volume 29., ACM (1999) 251–262

[9] Fountoulakis, N., Friedrich, T., Hermelin, D.: On the average-case complexity of parameterized clique. arXiv:1410.6400v1 (2014)

[10] Fountoulakis, N., Friedrich, T., Hermelin, D.: On the average-case complexity of parameterized clique. Theoretical Computer Science **576** (apr 2015) 18–29

[11] Friedrich, T., Krohmer, A.: Cliques in hyperbolic random graphs. In: INFOCOM proceedings 2015, IEEE (2015) 1544–1552

[12] Friedrich, T., Krohmer, A.: Parameterized clique on inhomogeneous random graphs. Discrete Applied Mathematics **184** (mar 2015) 130–138

[13] Garey, M.R., Johnson, D.S., Garey, M.R.: Computers and Intractability: A Guide to the Theory of NP-Completeness. W H FREEMAN & CO (1979)

[14] Grochow, J.A., Kellis, M.: Network motif discovery using subgraph enumeration and symmetry-breaking. In: In RECOMB. (2007) 92–106

[15] Heydari, H., Taheri, S.M.: Distributed maximal independent set on inhomogeneous random graphs. In: 2017 2nd Conference on Swarm Intelligence and Evolutionary Computation (CSIEC), IEEE (mar 2017)

[16] van der Hofstad, R., van Leeuwaarden, J.S.H., Stegehuis, C.: Optimal subgraph structures in scale-free networks. arXiv:1709.03466 (2017)

[17] Janson, S., Łuczak, T., Norros, I.: Large cliques in a power-law random graph. Journal of Applied Probability **47**(04) (dec 2010) 1124–1135

[18] Jeong, H., Tombor, B., Albert, R., Oltvai, Z.N., Barabási, A.L.: The large-scale organization of metabolic networks. Nature **407**(6804) (2000) 651–654

[19] Karp, R.M.: Reducibility among combinatorial problems. In: Complexity of computer computations. Springer (1972) 85–103

[20] Kashtan, N., Itzkovitz, S., Milo, R., Alon, U.: Efficient sampling algorithm for estimating subgraph concentrations and detecting network motifs. Bioinformatics **20**(11) (2004) 1746–1758

[21] Leskovec, J., Krevl, A.: SNAP Datasets: Stanford large network dataset collection. `http://snap.stanford.edu/data` (2014) Date of access: 14/03/2017.

[22] Niu, X., Sun, X., Wang, H., Rong, S., Qi, G., Yu, Y.: Zhishi.me - weaving chinese linking open data. In: The Semantic Web – ISWC 2011. Springer Nature (2011) 205–220

[23] Norros, I., Reittu, H.: On a conditionally poissonian graph process. Adv. Appl. Probab. **38**(01) (2006) 59–75

[24] Omidi, S., Schreiber, F., Masoudi-Nejad, A.: MODA: An efficient algorithm for network motif discovery in biological networks. Genes & Genetic Systems **84**(5) (2009) 385–395

[25] Park, J., Newman, M.E.J.: Statistical mechanics of networks. Phys. Rev. E **70** (2004) 066117

[26] Schreiber, F., Schwobbermeyer, H.: MAVisto: a tool for the exploration of network motifs. Bioinformatics **21**(17) (jul 2005) 3572–3574

[27] Vázquez, A., Pastor-Satorras, R., Vespignani, A.: Large-scale topological and dynamical properties of the internet. Phys. Rev. E **65** (2002) 066130

[28] Williams, V.V., Wang, J.R., Williams, R., Yu, H.: Finding four-node subgraphs in triangle time. In: Proceedings of the Twenty-sixth Annual ACM-SIAM Symposium on Discrete Algorithms. SODA '15, Philadelphia, PA, USA, Society for Industrial and Applied Mathematics (2015) 1671–1680