

# A Hybrid Sampling Scheme for Triangle Counting

John Kallaugher

`jmgk@cs.utexas.edu`

The University of Texas at Austin

Eric Price

`ecprice@cs.utexas.edu`

The University of Texas at Austin

## Abstract

We study the problem of estimating the number of triangles in a graph stream. No streaming algorithm can get sublinear space on all graphs, so methods in this area bound the space in terms of parameters of the input graph such as the maximum number of triangles sharing a single edge. We give a sampling algorithm that is additionally parameterized by the maximum number of triangles sharing a single vertex. Our bound matches the best known turnstile results in all graphs, and gets better performance on simple graphs like  $G(n, p)$  or a set of independent triangles.

We complement the upper bound with a lower bound showing that no sampling algorithm can do better on those graphs by more than a log factor. In particular, any insertion stream algorithm must use  $\sqrt{T}$  space when all the triangles share a common vertex, and any sampling algorithm must take  $T^{1/3}$  samples when all the triangles are independent. We add another lower bound, also matching our algorithm’s performance, which applies to *all* graph classes. This lower bound covers “triangle-dependent” sampling algorithms, a subclass that includes our algorithm and all previous sampling algorithms for the problem.

Finally, we show how to generalize our algorithm to count arbitrary subgraphs of constant size.

# 1 Introduction

Estimating the number of triangles in a graph is a fundamental graph problem. It is the smallest graph structure that cannot be counted directly from local information at each node, so it is a useful benchmark for understanding the power of various models of graph computation. At the same time, the number of triangles is itself a useful piece of information, and one of the standard features for understanding networks [ML12].

We consider the problem of triangle counting on graphs, under three different models of streaming access to the graph. In increasing order of restrictiveness, these are: *insertion-only* streams, where the edges  $E$  of  $G$  are added in adversarial order but are not deleted; *turnstile* streams, where edges may be inserted and deleted, and is equivalent (with some restrictions) to linear sketches [LNW14]; and nonadaptive *sampling*, where the algorithm specifies a distribution over sets  $S$  and receives  $E \cap S$ , for which it pays a “sample complexity” of  $\mathbb{E}|E \cap S|$ . Sampling is a more restrictive model than turnstile, because it can be implemented using linear sketches by composing an  $O(|E \cap S|) \times E$  noiseless sparse recovery sketch with a diagonal  $E \times E$  subsampling matrix. At the same time, many turnstile algorithms for graphs can be cast as sampling algorithms, and sampling algorithms have some advantages over general linear sketches, such as the ability to filter the graph after receiving the sketch, so we find it useful to distinguish the two models.

In the triangle estimation problem, one observes a graph  $G$  with  $T$  triangles and would like to output  $\bar{T}$  such that  $|T - \bar{T}| \leq \epsilon T$  with probability  $1 - \delta$ , using as little space/sample complexity as possible.

In the following discussion, we suppose that algorithms using a parameter know the optimal value of that parameter up to constant factors. We also use  $O^*(f)$  to hide factors of  $\epsilon$ ,  $\log \frac{1}{\delta}$ , and  $\log n$ .

**Streaming algorithms for triangle counting.** The problem of estimating triangles from a graph stream was introduced in [BKS02], which gave an  $O^*\left(\left(\frac{mn}{T}\right)^3\right)$  space algorithm based on estimating frequency moments in the insertion-only model. This was improved in [BFL<sup>+</sup>06] to  $O^*\left(\frac{mn}{T}\right)$  in the same model. An incomparable algorithm was given in [JG05] that uses  $O^*\left(\frac{md^2}{T}\right)$  space in the insertion-only model, where  $d$  is the maximum degree, based on subsampling the edges at rate  $d/T$  and storing all subsequent edges that touch the sampled edges.

None of the above algorithms are sublinear in  $m$  for all graphs. Unfortunately, as shown in [BOV13], this is inherent: even in the insertion model, no streaming algorithm can distinguish 0 triangles from  $T$  triangles in  $o(m)$  space for all graphs. The hard case involves a graph where all the triangles use a single edge. Since graphs of interest typically don’t involve that structure, a natural question arose of finding sublinear algorithms when  $\Delta_E$ , the maximum number of triangles sharing a common edge, is  $o(T)$ .

The first algorithms to achieve this, [TKMF09, TKM11], work by subsampling the edges with uniform probability  $p$ . The expected number of triangles is  $p^3 T$ , so one needs  $p > \frac{\Delta_E}{T} + \frac{1}{T^{1/3}}$  in order to (a) sample heavy edges with good probability and (b) sample at least one triangle in expectation. In fact this is sufficient, giving an algorithm with space  $O^*\left(m \left(\frac{\Delta_E}{T} + \frac{1}{T^{1/3}}\right)\right)$ .

In order to improve this, one would like a different sampling scheme that increases the chance of sampling edges in a triangle at the same time. In [PT12], this is done by randomly coloring the vertices of the graph with  $1/p$  colors, and keeping the monochromatic edges. This now samples  $p^2 T$  triangles in expectation, and indeed they show that  $O^*\left(m \left(\frac{\Delta_E}{T} + \frac{1}{T^{1/2}}\right)\right)$  samples suffice.

Paper	Space used	Model
[BKS02]	$(\frac{mn}{T})^3$	Insertion
[BFL <sup>+</sup> 06]	$\frac{mn}{T}$	Insertion
[JG05]	$\frac{md^2}{T}$	Insertion
[KP14, BFKP14]	$\frac{\sqrt{m}}{\alpha}$	Turnstile
[TKMF09, TKM11]	$m \left( \frac{\Delta_E}{T} + \frac{1}{T^{1/3}} \right)$	Sampling
[PT12]	$m \left( \frac{\Delta_E}{T} + \frac{1}{\sqrt{T}} \right)$	Sampling
This work	$m \left( \frac{\Delta_E}{T} + \frac{\sqrt{\Delta_V}}{T} + \frac{1}{T^{2/3}} \right)$	Sampling

Figure 1: Algorithms for triangle estimation. For simplicity, we drop dependencies on  $\epsilon$ ,  $\delta$ , and logarithmic factors. Here,  $\alpha$  denotes the transitivity coefficient,  $d$  denotes the maximum degree of any vertex, and  $\Delta_E/\Delta_V$  denote the maximum number of triangles sharing a common edge/vertex.

**Our results.** First, we show a distribution on graphs with  $\Delta_E = 1$  for which  $\Omega(\frac{m}{\sqrt{T}})$  space is required to estimate the number of triangles to within 25%, even in the insertion-only model. This shows that the [PT12] bound is tight in general. However, just as the  $\Omega(m)$  lower bound in [BOV13] involved unrealistic graphs where all triangles shared a single edge, our  $\Omega(m/\sqrt{T})$  lower bound involves unrealistic graphs where all triangles share a single vertex. So we consider algorithms additionally parameterized by  $\Delta_V$ , the maximum number of triangles sharing a common vertex.

We give a sampling algorithm using expected  $O^*(m \left( \frac{\Delta_E}{T} + \frac{\sqrt{\Delta_V}}{T} + \frac{1}{T^{2/3}} \right))$  samples.

**Theorem 1** (Triangle estimation). *If  $\tilde{T} \leq T$ , our algorithm obtains an  $(\epsilon, \delta)$  approximation to  $T$  while keeping  $O\left(\frac{m \log \frac{1}{\delta}}{\epsilon^2} \left(\frac{1}{\tilde{T}^{2/3}} + \frac{\sqrt{\Delta_V} \log \tilde{T}}{\tilde{T}} + \frac{\Delta_E \log \tilde{T}}{\tilde{T}}\right)\right)$  edges. If  $\tilde{T} > T$ , the algorithm either determines that  $\tilde{T} > T$  or obtains a  $(1 \pm \epsilon)$  approximation to  $T$  with probability  $1 - \delta$ .*

**Illustrative examples.** To compare our result with those in the literature, in Figure 2 we specialize the bounds to some illustrative example graphs. The *heavy edge* example is the lower bound from [BOV13] (and similar to one in [BFKP14]), where all triangles use a single edge. Any streaming algorithm requires  $\Omega(m)$  space in examples like this one, demonstrating the need for further parameterization, and most streaming algorithms match the bound. The *hub* example has a single vertex involved in  $n$  disjoint triangles. Here, we demonstrate that the  $\sqrt{n}$  achieved by [PT12] was optimal even for insertion-only streams.

We then consider graphs where most triangles do not overlap. One natural example is the Erdős-Rényi random graph  $G(n, p)$ , with  $1 \geq p \gg \frac{1}{n}$  so the number of triangles is well concentrated. Here, we use  $O^*(1/p)$  samples, and show that this is optimal for any sampling method. Previous algorithms in the turnstile/sampling models took  $\Omega(\sqrt{n})$  space for dense graphs, and even in the insertion-only model they took  $\Omega(1/p^2)$  space; our bound implies  $O^*(1/p)$ . The next example we consider is a collection of  $n$  independent triangles. Here, our algorithm takes  $O^*(n^{1/3})$  samples, which we show is optimal for any sampling algorithm, and improves upon the previous bound of  $\sqrt{n}$ . It is an interesting question whether  $n^{1/3}$  space is necessary in less restrictive streaming models.





Method	Model	Heavy edge 	Hub 	$G(n, p)$ 	Independent 
[BKS02]	Insertion	$n^3$	$n^3$	$1/p^6$	$n^3$
[BFL <sup>+</sup> 06]	Insertion	<b><math>n</math></b>	$n$	$1/p^2$	$n$
[JG05]	Insertion	$n^2$	$n^2$	$n$	<b>1</b>
[KP14, BFKP14]	Turnstile	$n^{3/2}$	$\sqrt{n}$	$n/\sqrt{p}$	$\sqrt{n}$
[TKMF09, TKM11]	Sampling	<b><math>n</math></b>	$n^{2/3}$	$n$	$n^{2/3}$
[PT12]	Sampling	<b><math>n</math></b>	$\sqrt{n}$	$\sqrt{n/p}$	$\sqrt{n}$
This paper	Sampling	<b><math>n</math></b>	$\sqrt{n}$	<b><math>1/p</math></b>	<b><math>n^{1/3}</math></b>
Lower bounds	Insertion	$n$ [BOV13, BFKP14]	$\sqrt{n}$	?	1
	Sampling	same	same	<b><math>1/p</math></b>	<b><math>n^{1/3}</math></b>

Figure 2: Results on specific graphs, ignoring logarithmic factors. For upper bounds, entries in bold are optimal in the computational model of the row. For lower bounds, entries in bold are new. The lower bound instances are subsets of the illustrated graphs containing a constant fraction of the triangles.

**Instance lower bounds.** It takes some care to define lower bounds for instances in a meaningful way. Let us consider trying to solve triangle counting for a specific class of graphs  $G_X$ , where  $X$  is some set of parameters (such the number of triangles in the hub graph). We would like to avoid “cheating” algorithms; for instance, a hub graph with  $m$  edges has  $m/3$  triangles, so one could estimate the number of triangles in the class of hub graphs by simply counting the edges, which takes 1 word in the turnstile model.

One attempt to avoid this would be to show the difficulty of distinguishing each graph  $G \in G_X$  from an alternative graph  $G'$  with a significantly different number of triangles. This is too weak:  $G'$  can introduce a difficult subproblem that is unnatural for the class  $G_X$ , so the lower bound doesn't really represent the difficulty of  $G_X$ . For example, when  $G$  has  $n$  independent triangles, we can have  $G'$  be  $G$  plus a clique on  $n^{1/3}$  vertices. A sampling algorithm using less than  $n^{1/3}$  edges would entirely miss the clique and be unable to distinguish the two, giving a lower bound by this definition, but not a satisfactory one: the alternative hard instance  $G'$  doesn't have the independent-triangle structural property we expect. To preserve the structure of the graph class, we would like to only consider graphs  $G'$  that are subgraphs of  $G$ .

The broader question here is, given that existing algorithms out-perform the  $\Omega(m)$  lower bound by adding the extra parameter  $\Delta_E$ , and we in turn out-perform these by adding the extra parameter  $\Delta_V$ , where should we draw the line? What are the correct set of parameters?

We suggest that an algorithm can correctly count triangles for a graph  $G$  with  $S$  samples and error  $\epsilon T(G)$ , it should be able to correctly count triangles for any subgraph  $G'$  of  $G$  with error  $\epsilon T(G)$ . (so additive error should be preserved, but not necessarily multiplicative error) Alternatively stated, the only thing that makes counting triangles in a larger graph easier is the fact that the larger graph may have more triangles, and thus have more tolerance for (additive) error.

We bolster this intuition by observing that existing algorithms all depend on two sets of param-

eters:  $T$  itself, in which their complexity is decreasing, and a set of monotonic<sup>1</sup> graph functions, in which their complexity is increasing (e.g.  $\Delta_V$  or  $\Delta_E$ )<sup>2</sup>.

We also assume that a triangle counting algorithm should be resilient to vertex labels being arbitrarily permuted, as it should not depend on knowing beforehand that certain vertices are “special.” This gives us our definition of “solving” an instance.

**Definition 2.** *Let  $G$  be a graph. We say an algorithm solves  $G$  with  $S$  space/samples and  $(\epsilon, \delta)$  error if, for any  $G'$  isomorphic to some subgraph of  $G$ , the algorithm returns  $T(G') \pm \epsilon T(G)$  with  $1 - \delta$  probability, using no more than  $S$  space/samples.*

We now define the *instance-optimum* for the space/sample complexity of solving  $G$ .

**Definition 3.** *For any given streaming model,  $\text{INSTOPT}(G, \epsilon, \delta)$  is the least amount of samples/space such that some algorithm solves  $G$  with  $\text{INSTOPT}(G, \epsilon, \delta)$  samples/space.*

An instance lower bound, then, is a lower bound on  $\text{INSTOPT}(G_X, \epsilon, \delta)$ . The lower bounds in Figure 2 use this definition.

It is an interesting question whether  $\Omega(n^{1/3})$  space is necessary for  $n$  independent triangles under turnstile or insertion streams. Our hard instance consists of randomly coloring the vertices with two colors, and either choosing the monochromatic or dichromatic edges. We do not know how to solve this instance with a turnstile streaming algorithm, but reductions from communication complexity seem to want three-party communication lower bounds, which are difficult to show. One can easily solve this particular instance in insertion streams, but similar instances of independent triangles with extra edges may be hard.

**Instance-optimal lower bound.** The results in Figure 2 show that our algorithm performs well on several natural graphs, but what about other graphs? Is there a more refined parameterization that would again yield large improvements on another class of graphs?

With some caveats, we show that our algorithm is optimal for *all* graphs. We need to refine our parameterization slightly, because currently the bound for a graph containing  $(1 - \epsilon/2)T$  independent triangles and a heavy edge with  $\epsilon T/2$  triangles would depend on the heavy edge, even though an upper bound could just skip it and remain within  $1 \pm \epsilon$  accuracy. We therefore define  $\Delta_{E,\epsilon}$  to be like  $\Delta_E$  but with the  $\epsilon T$  triangles contributing the most to  $\Delta_E$  removed, and  $\Delta_{V,\epsilon}$  similarly (for a precise definition, see Definition 16). This turns out to be a sufficient parameterization.

We show that any *triangle dependent* sampling algorithm—one that depends only on the set of triangles it samples—must use a set of samples with the expected dependence on  $\Delta_{E,2\epsilon}$  and  $\Delta_{V,2\epsilon}$ . All existing sampling algorithms for triangle counting are triangle dependent, but we cannot rule out better non-triangle-dependent algorithms.

**Theorem 4** (Triangle-dependent sampling bound). *For any constant  $\epsilon$  and for any graph  $G$ ,*

$$\text{INSTOPT}(G, \epsilon, 1/10) = \Omega \left( m \left( \frac{1}{T^{2/3}} + \frac{\sqrt{\Delta_{V,2\epsilon}}}{T} + \frac{\Delta_{E,2\epsilon}}{T} \right) \right)$$

*in the setting of triangle-dependent sampling algorithms.*

We then show that our upper bound can depend on  $\Delta_{E,\epsilon/24}$  and  $\Delta_{V,\epsilon/24}$ . Therefore, for any constant  $\epsilon$ , our algorithm calculates an  $(\epsilon, 1/10)$  approximation to the triangle count with  $O^*(\text{INSTOPT}(G, \epsilon/48, 1/10))$  samples.

<sup>1</sup>Meaning, in this context, that if  $A$  is a subgraph of  $B$ ,  $f(A) \leq f(B)$

<sup>2</sup>The only apparent exception is those algorithms which depend on the transitivity coefficient  $\alpha$ , but it holds if we replace  $\alpha$  with  $T/P_2$ ,  $P_2$  being the number of wedges in the graph.

**Beyond triangles.** We also give a generalization of our upper bound to counting arbitrary subgraphs of constant size. We give an algorithm to estimate  $M$ , the number of instances of a fixed size- $s$  subgraph, from a sample of the edges.

**Theorem 5** (Subgraph estimation). *Let  $f_\ell$  be the fraction of pairs of subgraphs that intersect at  $\ell$  vertices. We show how to find a  $1 + \epsilon$  factor approximation to  $M$  with probability  $1 - \delta$ , using order*

$$m \frac{\log(1/\delta)}{\epsilon^2} \log M \left( \sum_{\ell=2}^s f_\ell^{2/\ell} + f_\ell^{\frac{1}{\ell-1}} f_1^{1-\frac{1}{\ell-1}} \right)$$

*samples in expectation.*

For comparison, simple vertex sampling would replace the sum with  $\sum_{\ell=1}^s f_\ell^{2/\ell}$ ; our bound is always better. In the context of triangles, this difference is why we get  $\frac{\sqrt{\Delta_V}}{T}$  rather than  $\left(\frac{\Delta_V}{T}\right)^2$ , which was important in the hub case.

We get this bound using a similar scheme to our triangle estimation algorithm, removing a direct dependence on  $f_1$  by treating “heavy” vertices specially. In the case of  $s = 3$ , this is equivalent to our theorem up to a log factor.

**Other related work.** Another line of work parameterizes the space complexity using the *transitivity coefficient*  $\alpha$ , defined as the fraction of wedges that are completed into triangles. In [KP14, BFKP14] it is shown how to get  $O^*\left(\frac{\sqrt{m}}{\alpha}\right)$  space in the insertion model, for graphs without isolated edges. In [JSP13] it was shown that  $\tilde{O}\left(\frac{m}{\epsilon^2 \sqrt{T}}\right)$  space suffices in insertion streams to learn  $\alpha$  to  $\pm \epsilon$ . In fact, as we note in Appendix I, both bounds for triangle counting are directly implied by the  $O^*\left(m\left(\frac{\Delta_E}{T} + \frac{1}{\sqrt{T}}\right)\right)$  bound of [PT12]. Since our bound improves upon [PT12], it also implies these bounds.

[CJ14] shows multipass algorithms take  $\tilde{\Theta}(m/\sqrt{T})$  space for arbitrary graphs, giving an algorithm for two passes and a lower bound for a constant number of passes. [KMPT12] shows a three pass streaming algorithm using  $O(\sqrt{m} + m^{3/2}/T)$  space.

[ELRS15] considered the problem of triangle counting with query access to a graph. Similar to our algorithm, a simpler algorithm is modified to handle the impact of many vertices intersecting at a single triangle on the variance. The main difference is that, in [ELRS15], these “heavy” vertices are discarded without damaging the accuracy of the estimate, whereas we spend the bulk of our effort on attempting to estimate the number of triangles intersecting at each “heavy” vertex.

**Running time.** An  $O(m^{3/2})$  time algorithm was given in [IR78] to list all the triangles in a graph. This was improved for graphs with small arboricity by [CN85]. For *counting* triangles, [AYZ97] gave a different algorithm that improves the time to  $m^{\omega/(1+\omega)} \approx m^{1.41}$  using matrix multiplication. In [BPWZ14], it was shown how to extend this to *listing* triangles in  $o(m^{1.5})$  time when  $T = o(m^{3/2})$ . Other works, such as [SW05, Lat08], have given combinatorial  $O(m^{3/2})$  time triangle listing algorithms that are more efficient in practice. Our algorithm’s running time is dominated by listing triangles in the subsampled graph, which we can do either using one of the  $O(m^{3/2})$  time algorithms or (in some cases) slightly faster via [BPWZ14]. Because our algorithm improves upon the number of edges necessary to approximate the triangle count, it also implies a faster method for approximately counting the number of triangles in a given graph (as in, e.g., [KMPT12]).

## 2 Overview of Techniques

### 2.1 Triangle Counting Algorithm

This algorithm is a modification of simple vertex sampling, where we sample each vertex in the graph with probability  $\frac{1}{\sqrt{k}}$ , and keep any edge between two sampled vertices. The estimate of  $T$  is  $k^{\frac{3}{2}}$  times the number of triangles sampled, and the sample complexity is  $m/k$ . This is appealing, because as we show in our proof of Theorem 13, *any* algorithm that samples edges at rate  $1/k$  has a constant chance of sampling less than  $\frac{T}{k^{3/2}}$  triangles.

In order to approximate the number of triangles well, we need an estimator with variance  $O(T^2)$ . The vertex sampling estimator has variance bounded by  $k^{\frac{3}{2}}T + k \sum_e T_e^2 + \sqrt{k} \sum_v T_v^2$ , by the fact that a pair of triangles intersecting at  $l$  vertices is  $k^{\frac{l}{2}}$  times more likely to be sampled than a pair of triangles sampled independently. Choosing  $k$  small enough for this to be  $O(T^2)$  gives sample rate

$$\frac{1}{k} \approx \frac{1}{T^{2/3}} + \frac{\Delta_E}{T} + \left(\frac{\Delta_V}{T}\right)^2$$

The first and second terms are optimal, as we see in the independent triangles graph and heavy edge graph, but the third can be improved.

This is the term that makes vertex sampling fail on the hub case, when there is a single vertex  $v$  s.t.  $T_v = T$ . The vertex sampling algorithm will consistently fail to estimate the triangle count accurately in this case, as  $v$  will only be sampled with probability  $1/\sqrt{k}$ , and so usually we will miss all the triangles, and occasionally we will overestimate the triangle count by a factor of  $\sqrt{k}$ . More generally, our issue is vertices where  $T_v$  is large, and in particular vertices where it is larger than  $\frac{T}{\sqrt{k}}$ , as the total contribution to the variance of vertices with  $T_v$  smaller than this is  $O(T^2)$ .

We can deal with the hub case by extended the sampling in the following way. After sampling vertices, in addition to taking all the edges between sampled vertices, we can also take  $1/\sqrt{k}$  of the edges between sampled vertices and unsampled vertices. Now, even when the central vertex  $v$  is not sampled, each triangle in the hub has a  $1/k^2$  chance of being sampled (if both other vertices, and both edges between  $v$  and those vertices, are picked). This is independent for the different triangles in the hub, so we will find a triangle when  $k \approx \sqrt{T}$ . The resulting  $m/\sqrt{T}$  sample complexity is optimal by our lower bound.

So we could handle our vertices one of two ways—for our “light” vertices, we will use vertex sampling, and for our “heavy” ( $T_v \geq \frac{T}{\sqrt{k}}$ ) vertices, we will use the scheme above. In order to identify the lightest of the heavy vertices, as we are sampling their triangles at rate  $k^{-2}$ , we would need  $k^2 = \frac{T}{\sqrt{k}}$  and so  $k < T^{\frac{2}{5}}$ . Can we do better?

We can think of vertex sampling and the “hub scheme” as two ends of a continuum. For a given “weight”  $x$ , we can take two samples  $S_1$  and  $S_2$  of  $V$  where each vertex appears in  $S_1$  with probability  $1/\sqrt{k}$  and appears in  $S_2$  with probability  $x/\sqrt{k}$ , then sample each edge in  $S_1 \times S_2$  with probability  $1/x$ . The  $x = 1$  case is vertex sampling, and the  $x = \sqrt{k}$  case is our hub scheme. The trade-off is that small  $x$  makes it possible to completely miss important vertices, but high  $x$  means that we sample fewer triangles overall.

Consider the “many-hubs” case where we have  $T/\Delta_V$  hubs involved in  $\Delta_V$  triangles each. We need  $x \geq \sqrt{k}\Delta_V/T$  to sample them reliably in the weight  $x$  scheme, and we will get  $\frac{T}{k^{3/2}x}$  triangles in expectation. When  $x$  is minimized according to the first constraint, this is  $\frac{T^2}{\Delta_V k^2}$  triangles in expectation. For the variance to be small, we need this to be at least a constant, as happens for  $k \leq \frac{\sqrt{\Delta_V}}{T}$ . This explains the  $\frac{\sqrt{\Delta_V}}{T}$  bound in our algorithm.

The above discussion applies when  $x$  is optimized for a given graph. Our full algorithm runs the scheme for all  $\log k$  different scales of  $x$  and combines the results. This lets us improve the variance bound to order  $T^2 + k \sum_e T_e^2 + k^2 \frac{\sum_v T_v^2}{T} 3$  at the cost of sampling each edge with probability  $\frac{\log k}{k}$  instead of  $\frac{1}{k}$ . In particular, we improve our performance on a hub graph from sampling at rate 1 to rate  $\frac{\log T}{\sqrt{T}}$ , while on graphs with less triangle-heavy vertices, we can achieve the bound of a  $\frac{1}{T^{3/2}}$  sampling rate.

As we only want to use those sampling rates appropriate to the scales of vertex actually present in the graph (as our performance depends on the greatest scale), we parameterize our algorithm by  $\omega$ , the minimum weight we will put on a vertex, which is  $\min \left\{ \frac{T^2}{\sum_v T_v^2}, \sqrt{k} \right\}$ .

Our algorithm conceptually runs as in two passes over the stream. First, for each vertex  $v$  of the graph, we calculate an estimate  $\mathcal{T}_v \in \left\{ 0, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil \right\} \cup \{\text{LIGHT}\}$  of  $\left\lceil \log \frac{T}{\omega T_v} \right\rceil$ , with  $\mathcal{T}_v = \text{LIGHT}$  when  $T_v$  is believed to be  $< \frac{\tilde{T}}{\sqrt{k}}$ , and 0 when it is believed to be  $> \frac{T}{\omega}$ . In the second pass, we estimate  $T_L$ , the number of triangles  $t = (u, v, w)$  s.t.  $\mathcal{T}_u = \mathcal{T}_v = \mathcal{T}_w = \text{LIGHT}$ , and  $T_H$ , the number of triangles using at least one vertex  $v$  s.t.  $\mathcal{T}_v \neq \text{LIGHT}$ .

We will show that it is possible to perform both conceptual passes in one pass over the data, by only calculating those  $\mathcal{T}_v$  which are needed for the second pass.

## 2.2 Instance Lower Bounds

We recall our definition of *instance-optimum* for a class of graphs.

**Definition 2.** Let  $G$  be a graph. We say an algorithm solves  $G$  with  $S$  space/samples and  $(\epsilon, \delta)$  error if, for any  $G'$  isomorphic to some subgraph of  $G$ , the algorithm returns  $T(G') \pm \epsilon T(G)$  with  $1 - \delta$  probability, using no more than  $S$  space/samples.

**Definition 3.** For any given streaming model,  $\text{INSTOPT}(G, \epsilon, \delta)$  is the least amount of samples/space such that some algorithm solves  $G$  with  $\text{INSTOPT}(G, \epsilon, \delta)$  samples/space.

We demonstrate that it is sufficient to show that an algorithm cannot distinguish between two distributions on subgraphs of  $G$  with triangle counts separated by  $\Omega(T)$ .

**Definition 6** (Distinguishing). We say that an algorithm  $\mathcal{A}$  can distinguish two random graph distributions  $\mathcal{G}_1$  and  $\mathcal{G}_2$  if there exists  $f$  such that, for a pair of draws  $G_1$  and  $G_2$  from these distributions, and any relabelling of the vertices of  $G_1$  and  $G_2$ ,  $\Pr[f(\mathcal{A}(G_1)) = 1] \geq 3/4$  and  $\Pr[f(\mathcal{A}(G_2)) \neq 1] \geq 3/4$ .

**Lemma 7.** Let  $\mathcal{A}$  be an algorithm that solves triangle counting for a graph  $G$  with  $S$  space/samples and  $(\epsilon, 1/10)$  error. Then, for any two distributions  $\mathcal{G}_1, \mathcal{G}_2$  on subgraphs  $G_1$  and  $G_2$  of  $G$ , and  $C$  such that  $T(G_1) > C + \epsilon T(G)$  with  $\frac{9}{10}$  probability and  $T(G_2) < C - \epsilon T(G)$  with  $\frac{9}{10}$  probability,  $\mathcal{A}$  can distinguish them.

### 2.2.1 Heavy Edges Graph

**Definition 8** (Heavy Edges Graph). The heavy edges graph  $D_{r,d}$  consists of  $r$  copies of the following graph:  $d$  disjoint edges  $\{u_{2i}u_{2i+1}\}_{i=0}^{d-1}$ , one of which has both ends connected to a further  $d$  vertices  $\{v_i\}_{i=0}^{d-1}$ .

---

<sup>3</sup>Note that  $\sum_v T_v^2 \leq \Delta_v T$ .



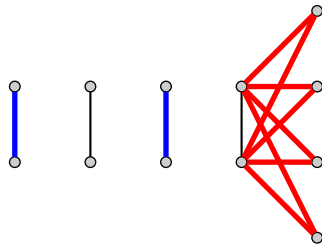
**Theorem 9** (Heavy Edges Lower Bound). *For sufficiently small constant  $\epsilon$ ,  $\text{INSTOPT}(D_{r,d}, \epsilon, 1/10) = \Omega(d)$  bits in the insertion-only model.*

We use a reduction shown in [BOV13]. We reduce from the indexing problem  $\text{Index}_n$  to the problem of distinguishing an  $rd$ -triangle subset of  $D_{r,d}$  and a 0-triangle subset.

$\text{Index}_n$  is defined as follows: Alice has a binary vector  $w$  of length  $n$ , and Bob has an index  $x \in [n]$ . Alice must send a message to Bob such that Bob can determine  $w[x]$ . By [CCKM10], the randomized communication complexity of this problem is  $\Omega(n)$ .

Alice can encode  $w$  in a subgraph of  $D_{1,d}$  by, for each  $i$  in  $\{0, \dots, n-1\}$ , including  $u_{2i}u_{2i+1}$  iff  $w_i = 1$ . Bob then adds the  $d$  vertices  $\{v_i\}_{i=0}^{d-1}$  to the graph, connecting each of them to  $u_{2x}$  and  $u_{2x+1}$ . The resulting graph will have  $d$  triangles if  $w_x = 1$  and 0 otherwise.

This result extends to arbitrary  $r$  by letting Alice and Bob each repeat their encoding  $r$  times.



Encoding the string 1010 in  $D_{1,4}$ , Alice sends the edges in blue, connecting the  $i^{\text{th}}$  pair of vertices if the  $i^{\text{th}}$  bit of the string is 1. Bob then queries the 4<sup>th</sup> position by adding the edges in red, connecting the  $i^{\text{th}}$  pair to  $d = 4$  wedges. The graph will contain  $d$  triangles if the  $i^{\text{th}}$  bit is 1, and 0 otherwise.

Figure 3: Heavy edge graph, showing that  $\Omega(m \frac{\Delta E}{T})$  is necessary for insertion streams.

## 2.2.2 Hubs Graph

**Definition 10** (Hubs Graph). *The hubs graph  $H_{r,d}$  consists of  $r$  copies of the following graph: a single vertex which participates in  $d$  edge-disjoint triangles.*

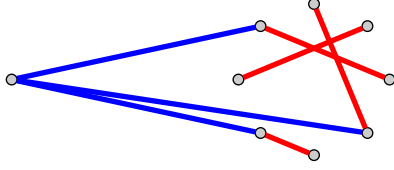
**Theorem 11** (Hubs Lower Bound). *For sufficiently small constant  $\epsilon$ ,  $\text{INSTOPT}(H_{r,d}, \epsilon, \delta) = \Omega(\sqrt{d})$  bits in the insertion-only model.*

For a single hub, we consider the following communication problem: there are  $d$  people who might go on a trip. Alice knows who is going on the trip, and Bob knows who among them are couples. Alice must send a message to Bob that lets him determine whether the trip is at least  $2/3$  couples, or at most  $1/3$  couples. Intuitively, this should require  $\Omega(\sqrt{d})$  communication, as if Bob learns fewer than  $\sqrt{d}$  of the people who are going on the trip, he is likely to not learn of any of the couples. We can then consider the hub edges as the identities of people who are going on the trip, and the other edges as the identities of the couples.

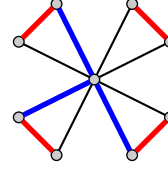
Formally, we use a reduction from the variant of the Boolean Hidden Matching problem  $\text{BHM}_n$  set out in [KR06], which has  $\Omega(\sqrt{n})$  randomized communication complexity.

## 2.2.3 Independent Triangles

**Definition 12** (Independent Triangles Graph). *The independent triangles graph  $I_n$  consists of  $n$  vertex-disjoint triangles.*



(a) Encoding a “trip”.



(b) The same encoding, now shown as a subgraph of  $H_{1,4}$ .

We have one hub vertex, and a set of spoke vertices, each corresponding to a potential traveller. Alice sends the edges in blue, drawing an edge from the hub vertex to each traveller on the trip. Bob adds the edges in red, drawing an edge between the members of each couple. The number of triangles in the resulting graph will be the same as the number of couples on the trip.

Figure 4: Hub graph, showing that  $\Omega(m \frac{\sqrt{\Delta V}}{T})$  is necessary for insertion streams.

**Theorem 13** (Independent Triangles Lower Bound). *For sufficiently small constant  $\epsilon$ ,  $\text{INSTOPT}(I_n, \epsilon, 1/10) = \Omega\left(n^{\frac{1}{3}}\right)$  samples in the sampling model.*

We prove this by showing that we need to sample edges with probability  $\Omega\left(\frac{1}{n^{2/3}}\right)$  to achieve a constant chance of sampling *any* triangle.

Intuitively, it seems that this should be sufficient—a sampling algorithm which fails to sample any triangles should not be able to count the number of triangles in a graph. However, we have been unable to prove this in the general case, as for arbitrary graphs, finding two subgraphs with a large enough separation in triangle count that cannot be distinguished without sampling any triangles turns out to be difficult. For instance, the pair would need to have the same number of edges, as otherwise they could be distinguished simply by edge counting, and the same number of wedges, as otherwise they could be distinguished by an algorithm that samples wedges but not triangles.

In this specific case, however, we can show that there are two different (distributions on) subgraphs which satisfy this requirement.

Conditioned on failing to sample any triangles (and therefore, in  $I_n$ , any cycles), we show that, for a random 2-coloring  $\chi$ , the distributions of the following two graphs under the sampling scheme are identical:

$$I_1 = (V(I_n), \{uv \in E(I_n) | \chi(u) = \chi(v)\})$$

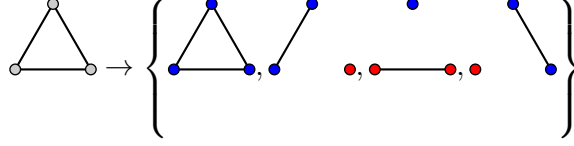
$$I_2 = (V(I_n), \{uv \in E(I_n) | \chi(u) \neq \chi(v)\})$$

Which, as  $I_1$  has a constant fraction of the triangles, and  $I_2$  has none of them, proves our theorem.

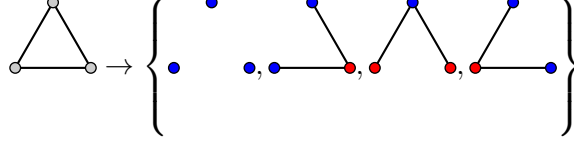
#### 2.2.4 $G_{n,p}$

**Theorem 14** ( $G_{n,p}$  Lower Bound). *There exists a constant  $C$  such that, provided  $p \geq \frac{C}{n}$ , and for sufficiently small constant  $\epsilon$ ,  $\text{INSTOPT}(G_{n,p}, \epsilon, 1/10) = \Omega\left(\frac{1}{p}\right)$  samples.*

We make use of the same pair of colorings as in the proof of the independent triangles lower bound, but as our triangles are no longer guaranteed not to intersect, we prove that we need  $\Omega\left(\frac{1}{p}\right)$  samples to sample *any* cycles from the graph. Subject to our sample being acyclic, it will be



(a) The four (equally likely) possible results for applying the first subgraphing scheme to a single triangle.



(b) The four (equally likely) possible results for applying the second subgraphing scheme to a single triangle.

A sampling scheme that only looks at one or two of the edges will observe identical distributions under the two regimes, so a sampling scheme that distinguishes the regimes must sample at least one complete triangle.

Figure 5: Independent triangles, showing that  $\Omega(m/T^{2/3})$  is necessary for sampling algorithms.

identically distributed between both colorings, despite their differing by a constant fraction of the graph's triangles, and so it follows that we will not be able to count triangles.

### 2.3 Instance-Optimality for Triangle-Dependent Sampling Algorithms

We now present a lower bound which applies to all graph instances, but only a subclass of sampling algorithms, specifically those which depend on actually sampling at least one triangle to find a  $(1 \pm \epsilon)$  estimate of the number of triangles in a graph.

**Definition 15.** Let  $\mathcal{A}$  be a sampling algorithm for counting triangles. We say that  $\mathcal{A}$  is a triangle-dependent sampling algorithm if, for all graphs  $G$ ,  $\mathcal{A}(G)$  depends only on the set of triangles sampled by  $\mathcal{A}$ .

Note that this definition encompasses all existing sampling algorithms for triangle counting.

Note that this definition encompasses all the strategies for counting triangles by edge sampling mentioned earlier, all of which depend on sampling edges by some strategy and then weighting the triangles sampled this way.

Our method is to show, that each of the three parameters in our graph are, in any graph, necessary for sampling triangles. This requires extending the definition of  $\Delta_E$  and  $\Delta_V$  to allow excluding  $\epsilon T$  of the “heaviest” vertices or edges. This is necessary because a graph may, for instance, have a single edge with  $\epsilon/2$  triangles, which contributes to the variance of the estimator, but which will not prevent the algorithm accurately estimating the triangle count if it is not sampled.

**Definition 16.** For any graph  $G$ ,  $\epsilon > 0$ , let the vertices  $v \in V(G)$  be ordered as  $(v_i)_{i \geq 0}$  in descending order of  $T_v$ , and the edges  $e \in E(G)$  be ordered as  $(e_i)_{i \geq 0}$  in descending order of  $T_e$ . Then let  $H_V$ ,  $H_E$  be the maximal prefixes of  $(v_i)_{i \geq 0}$ ,  $(e_i)_{i \geq 0}$  such that  $\sum_{v \in H_V} T_v$ ,  $\sum_{e \in H_E} T_e \leq \epsilon T$ . We define  $\Delta_{V,\epsilon}(G) = \max_{v \notin H_V} T_v$ ,  $\Delta_{E,\epsilon}(G) = \max_{e \notin H_E} T_e$ .

When the graph meant is unambiguous, we will omit the parameter  $G$ .

**Theorem 4** (Triangle-dependent sampling bound). For any constant  $\epsilon$  and for any graph  $G$ ,

$$\text{INSTOPT}(G, \epsilon, 1/10) = \Omega \left( m \left( \frac{1}{T^{2/3}} + \frac{\sqrt{\Delta_{V,2\epsilon}}}{T} + \frac{\Delta_{E,2\epsilon}}{T} \right) \right)$$

in the setting of triangle-dependent sampling algorithms.

Analyzing our algorithm in terms of the variance does not allow us to reach this bound. However, without altering the algorithm itself, we can refine the analysis by “cutting off” a small number (less than  $\epsilon T$  times a small constant) of the heaviest vertices and edges. This works because their contribution to the estimate  $\bar{T}$  is always positive, and so we may split  $\bar{T}$  into  $\bar{T}_\epsilon$  (representing  $\bar{T}$  less the contribution of these edges and vertices), which we bound by Chebyshev’s inequality as usual, and  $(\bar{T} - \bar{T}_\epsilon)$ , which we bound by Markov’s inequality. This allows us to replace the  $\Delta_V, \Delta_E$  terms in our bound with  $\Delta_{V,\epsilon}, \Delta_{E,\epsilon}$ , so that, for constant  $\epsilon$ , we can match the lower bound up to a log factor and by a constant<sup>4</sup> factor in  $\epsilon$ .

**Theorem 17** (Refined triangle estimation upper bound). *If  $\tilde{T} \leq T$ , and  $\epsilon > 0$ , our algorithm obtains an  $(\epsilon, \delta)$  approximation to  $T$  while keeping*

$$O\left(\frac{m \log \frac{1}{\delta}}{\epsilon^2} \left(\frac{1}{\tilde{T}^{2/3}} + \frac{\sqrt{\Delta_{V,\epsilon/24}} \log \tilde{T}}{\tilde{T}} + \frac{\Delta_{E,\epsilon/24} \log \tilde{T}}{\tilde{T}}\right)\right)$$

*edges. If  $\tilde{T} > T$ , the algorithm either determines that  $\tilde{T} > T$  or obtains a  $(1 \pm \epsilon)$  approximation to  $T$  with probability  $1 - \delta$ .*

## 2.4 Algorithm for General Constant-Size Subgraphs

We show that the algorithm in this paper can, with small modifications, be generalized to count the number of copies of an arbitrary constant-size subgraph  $A$  in  $G$ . As in the triangle case, we start with the algorithm in which every vertex is sampled with probability  $\frac{1}{\sqrt{k}}$ , and edges between sampled vertices are kept. This will give an estimator with variance:

$$O\left(\sum_{l=1}^s (C_l k^{\frac{l}{2}})\right)$$

Where we define  $C_i$  as follows: for any  $S \subseteq V(G)$ , we define  $\alpha(S)$  defined as the set of copies of  $A$  that use all the vertices in  $S$  and  $M_S$  as  $|\alpha(S)|$ .  $C_i$  is then defined as  $\sum_{|S|=i} M_S$ . Note that in the case where  $A$  is a triangle,  $C_1 = \sum_v T_v^2$ ,  $C_2 = \sum_e T_e^2$ , and  $C_3 = T$ .

As in the triangle case, we will eliminate the  $C_1$  term, by estimating the number of subgraphs involving  $v$  at each  $v$  and reducing the “weight” applied to  $v$  accordingly. In the triangle case, this increased the  $C_3$  (equivalently,  $T$ ) term—in the general case it will increase the  $C_l$  term for every  $l \geq 3$ . This is because, when we put less weight on a vertex, for any pair of subgraphs  $a_1, a_2$  which intersect at  $\geq 2$  edges adjacent to this vertex, the event of sampling  $a_1$  will become *more* correlated with the event of sampling  $a_2$ .

Consequently, our new variance will end up depending on  $C_1$  at every term of the sum, giving us variance:

$$M^2 + \sum_{l=2}^s C_l \left( k^{\frac{l}{2}} + k \left( \frac{C_1^+}{M^2} k \right)^{l-2} \right)$$

This then gives us our general subgraphs result.

<sup>4</sup>The distinction here is that multiplying  $\epsilon$  by a constant can cause a non-constant change in  $\Delta_{E,\epsilon}$

**Theorem 5** (Subgraph estimation). *Let  $f_\ell$  be the fraction of pairs of subgraphs that intersect at  $\ell$  vertices. We show how to find a  $1 + \epsilon$  factor approximation to  $M$  with probability  $1 - \delta$ , using order*

$$m \frac{\log(1/\delta)}{\epsilon^2} \log M \left( \sum_{\ell=2}^s f_\ell^{2/\ell} + f_\ell^{\frac{1}{\ell-1}} f_1^{1-\frac{1}{\ell-1}} \right)$$

*samples in expectation.*

## Acknowledgements

The authors would like to thank David Woodruff for a helpful pointer to the Boolean Hidden Matching problem [KR06].

## References

- [AYZ97] Noga Alon, Raphael Yuster, and Uri Zwick. Finding and counting given length cycles. *Algorithmica*, 17(3):209–223, 1997.
- [BFKP14] Laurent Bulteau, Vincent Froese, Konstantin Kutzkov, and Rasmus Pagh. Triangle counting in dynamic graph streams. *Algorithmica*, pages 1–20, 2014.
- [BFL<sup>+</sup>06] Luciana S Buriol, Gereon Frahling, Stefano Leonardi, Alberto Marchetti-Spaccamela, and Christian Sohler. Counting triangles in data streams. In *Proceedings of the twenty-fifth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 253–262. ACM, 2006.
- [BKS02] Ziv Bar-Yossef, Ravi Kumar, and D. Sivakumar. Reductions in streaming algorithms, with an application to counting triangles in graphs. In *Proceedings of the Thirteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '02*, pages 623–632, Philadelphia, PA, USA, 2002. Society for Industrial and Applied Mathematics.
- [BOV13] Vladimir Braverman, Rafail Ostrovsky, and Dan Vilenchik. How hard is counting triangles in the streaming model? In *Automata, Languages, and Programming*, pages 244–254. Springer, 2013.
- [BPWZ14] Andreas Björklund, Rasmus Pagh, Virginia Vassilevska Williams, and Uri Zwick. Listing triangles. In *Automata, Languages, and Programming*, pages 223–234. Springer, 2014.
- [CCKM10] Amit Chakrabarti, Graham Cormode, Ranganath Kondapally, and Andrew McGregor. Information cost tradeoffs for augmented index and streaming language recognition. In *Proceedings of the 51st FOCS*, pages 387–396. IEEE, 2010.
- [CJ14] Graham Cormode and Hossein Jowhari. A second look at counting triangles in graph streams. *Theoretical Computer Science*, 552:44–51, 2014.
- [CN85] Norishige Chiba and Takao Nishizeki. Arboricity and subgraph listing algorithms. *SIAM Journal on Computing*, 14(1):210–223, 1985.

- [eh14] emab (<http://math.stackexchange.com/users/74964/emab>). Number of triangles in a graph based on number of edges. Mathematics Stack Exchange, 2014. URL:<http://math.stackexchange.com/q/823650> (version: 2014-06-07).
- [ELRS15] Talya Eden, Amit Levi, Dana Ron, and C. Seshadhri. Approximately counting triangles in sublinear time. In *Proceedings of the 56th FOCS*, pages 614–633. IEEE, 2015.
- [IR78] Alon Itai and Michael Rodeh. Finding a minimum circuit in a graph. *SIAM Journal on Computing*, 7(4):413–423, 1978.
- [JG05] Hossein Jowhari and Mohammad Ghodsi. New streaming algorithms for counting triangles in graphs. In *Computing and Combinatorics*, pages 710–716. Springer, 2005.
- [JSP13] Madhav Jha, Comandur Seshadhri, and Ali Pinar. A space efficient streaming algorithm for triangle counting using the birthday paradox. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 589–597. ACM, 2013.
- [KMPT12] Mihail N Kolountzakis, Gary L Miller, Richard Peng, and Charalampos E Tsourakakis. Efficient triangle counting in large graphs via degree-based vertex partitioning. *Internet Mathematics*, 8(1-2):161–185, 2012.
- [KP14] Konstantin Kutzkov and Rasmus Pagh. Triangle counting in dynamic graph streams. In *Algorithm Theory–SWAT 2014*, pages 306–318. Springer, 2014.
- [KR06] I. Kerenidis and R. Raz. The one-way communication complexity of the Boolean Hidden Matching Problem. 2006.
- [Lat08] Matthieu Latapy. Main-memory triangle computations for very large (sparse (power-law)) graphs. *Theoretical Computer Science*, 407(1):458–473, 2008.
- [LNW14] Yi Li, Huy L Nguyen, and David P Woodruff. Turnstile streaming algorithms might as well be linear sketches. In *Proceedings of the 46th annual ACM Symposium on Theory of Computing*, pages 174–183. ACM, 2014.
- [ML12] Julian J McAuley and Jure Leskovec. Learning to discover social circles in ego networks. In *NIPS*, volume 2012, pages 548–56, 2012.
- [MV08] Jiří Matoušek and Jan Vondrák. The probabilistic method lecture notes, March 2008.
- [PT12] Rasmus Pagh and Charalampos E Tsourakakis. Colorful triangle counting and a mapreduce implementation. *Information Processing Letters*, 112(7):277–281, 2012.
- [SW05] Thomas Schank and Dorothea Wagner. Finding, counting and listing all triangles in large graphs, an experimental study. In *Experimental and Efficient Algorithms*, pages 606–609. Springer, 2005.
- [TKM11] Charalampos E Tsourakakis, Mihail N Kolountzakis, and Gary L Miller. Triangle sparsifiers. *J. Graph Algorithms Appl.*, 15(6):703–726, 2011.
- [TKMF09] Charalampos E Tsourakakis, U Kang, Gary L Miller, and Christos Faloutsos. Doulion: counting triangles in massive graphs with a coin. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 837–846. ACM, 2009.

# A Algorithm

## A.1 Definitions

Let  $G$  be a graph, where we receive  $E(G)$  as a stream of edges, with  $m = |E(G)|$ . Edges in  $E(G)$  will always be treated as undirected.

$T(G)$  is the number of triangles in  $G$ . For any  $v \in V(G)$ , let  $\tau(v)$  be the set of triangles in  $G$  involving  $v$ , and  $T_v = |\tau(v)|$ . For any  $e \in E(G)$ , let  $T_e$  be the number of triangles involving  $e$ . Then, let  $T_V(G) = \sum_v T_v^2$  and  $T_E(G) = \sum_e T_e^2$ . Where the graph  $G$  meant is unambiguous, we will omit the explicit parametrization by  $G$ . Note that  $T_V \leq T\Delta_V$  and  $T_E \leq T\Delta_E$ , respectively.

$k \in [0, 1]$  is our sampling parameter. We will show that the expected number of edges stored by the algorithm is  $O\left(\frac{m \log k}{k}\right)$ .

Let  $\tilde{T}$  be a proposed lower bound on  $T$ , and  $T_V^+$  an actual upper bound on  $T_V$ . If  $T_V \leq T_V^+$  and  $T \geq \tilde{T}$ , we want to be able to accurately estimate the number of triangles in the graph. If  $T_V \leq T_V^+$  and  $T < \tilde{T}$ , we want to be able to detect that  $T < \tilde{T}$ . Given these two parameters, we will define  $\omega := \min\left\{\frac{(\tilde{T})^2}{T_V^+}, \sqrt{k}\right\}$ .  $\omega$  will be the minimum “weight” we put on a vertex, and as we will use weights from  $\omega$  to  $\sqrt{k}$ , the total number of sampling schemes we use will be  $\log \frac{\sqrt{k}}{\omega}$ . This means that the bound on the expected number of edges stored by the algorithm can in fact be reduced to  $O\left(\frac{m}{k} \log \frac{\sqrt{k}}{\omega}\right)$ .

## A.2 First pass

### A.2.1 Outline

In our first pass over the graph, we attempt to estimate the “correct” weight to put on each vertex  $v$ ,  $\frac{T}{T_v}$ . As we do not know  $T$ , we use our lower bound  $\tilde{T}$ , which will allow us to bound the variance of our final estimate in terms of  $\tilde{T}$  and therefore  $T$ . As we do not have direct access to  $T_v$ , we look at the number of triangles that would be counted at  $v$  if it were assigned a given weight.

We consider  $\log \frac{\sqrt{k}}{\omega}$  possible weights that could be assigned to  $v$ . For each  $h \in \left\{0, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil\right\}$ , we consider the weight  $\omega 2^h$  ( $\omega$ , therefore, being the minimum weight we will assign any vertex). Let  $\mathcal{T}_v$  be our estimate of the “correct” choice of  $h$ . But what should this be? Letting  $X_v^{(h)}$  be the number of triangles counted at  $v$  when it is assigned weight  $2^h$ , the expected value of  $X_v^{(h)}$  will be  $\frac{2^{2h} T_v}{k^2}$ . So if we want  $\omega 2^h$  to be  $\frac{\tilde{T}}{T_v}$ , the expectation of  $X_v^{(h)}$  should be  $\frac{\omega \tilde{T} 2^h}{k^2}$ . We will therefore define  $\mathcal{T}_v$  to be the least  $h$  such that  $X_v^{(h)}$  is at least this high. If  $X_v^{(h)}$  is not this high for any  $h \in \left\{0, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil\right\}$ , we will conclude that  $X_v \leq \frac{\tilde{T}}{\sqrt{k}}$ , and so our algorithm should treat it as a “light” vertex, and so we set  $\mathcal{T}_v = \text{LIGHT}$ .

In order to attain the bounds we need on the distance of  $\mathcal{T}_v$  from  $\left\lceil \log \frac{T}{\omega T_v} \right\rceil$ , we will run this procedure twice, taking the lowest value of  $\mathcal{T}_v$ .

### A.2.2 Procedure

For  $i = 1, 2$ , and for each  $v \in V(G)$ , define  $\mathcal{T}_{v,i}$  as follows:

Let  $r_{D,i} : E \rightarrow [0, 1]$  be a uniformly random hash function.

Let  $d_i : V \rightarrow \{0, 1\}$  be a random hash function such that  $\forall v \in V$ :

$$d_i(v) = \begin{cases} 1 & \text{With probability } \frac{1}{\sqrt{k}}. \\ 0 & \text{Otherwise.} \end{cases}$$

For each  $v \in V(G)$ ,  $t = (u, v, w)$  a triangle in  $G$ , and  $h \in \{0, \dots, \lceil \log \frac{\sqrt{k}}{\omega} \rceil\}$ , define  $X_{v,i}^{(h,t)}$  to be 1 if:

$$\begin{aligned} d_i(u) &= 1 \\ d_i(w) &= 1 \\ r_{D,i}(vu) &< \frac{\omega 2^h}{\sqrt{k}} \\ r_{D,i}(vw) &< \frac{\omega 2^h}{\sqrt{k}} \end{aligned}$$

And 0 otherwise. Then let  $X_{v,i}^{(h)} = \sum_{t \in \tau(v)} X_{v,i}^{(h,t)}$ . We then define  $H_{v,i} = \left\{ h \in \left\{ 1, \dots, \lceil \log \frac{\sqrt{k}}{\omega} \rceil \right\} \mid X_{v,i}^{(h)} \geq \frac{\omega \tilde{T} 2^h}{k^2} \right\}$ .  $\mathcal{T}_{v,i}$  is then defined as follows:

$$\mathcal{T}_{v,i} = \begin{cases} \min H_{v,i} & \text{If } H_{v,i} \neq \emptyset. \\ \text{LIGHT} & \text{Otherwise.} \end{cases}$$

Then, for each  $v \in V(G)$ , we define  $\mathcal{T}_v$  to be LIGHT if  $\mathcal{T}_{v,1}$  and  $\mathcal{T}_{v,2}$  are LIGHT, and otherwise to be the smallest numerical value amongst  $\mathcal{T}_{v,1}, \mathcal{T}_{v,2}$ .

### A.3 Second Pass

#### A.3.1 Outline

We now use the scale estimates  $\mathcal{T}_v$  to determine our strategy for estimating  $T$ . If, for a vertex  $v$ ,  $\mathcal{T}_v = \text{LIGHT}$ , we believe that  $T_v \leq \frac{\tilde{T}}{\sqrt{k}}$ , so we use our naïve sampling method to estimate  $T_L$ , the number of triangles in  $G$  which use only such vertices. Otherwise, we sample  $v$  with probability  $2^{-\mathcal{T}_v}$ , and if we sample it, construct an estimate  $\overline{T}_v$  of  $T_v$  by assigning  $v$  weight  $2^{\mathcal{T}_v}$ . As this could lead to a triangle being counted up to three times (if all three of its vertices  $v$  had  $\mathcal{T}_v \neq \text{LIGHT}$ , it would be counted three times), we give each triangle either weight 1,  $\frac{2}{3}$ , or  $\frac{1}{3}$ , depending on how many of its vertices  $v$  have  $\mathcal{T}_v \neq \text{LIGHT}$ . This then lets us report our estimate  $\overline{T}_H$  of  $T_H$  as  $\sum_v \overline{T}_v$ .

#### A.3.2 Splitting the Graph

Let  $V_L = \{v \in V \mid \mathcal{T}_v = \text{LIGHT}\}$ , and let  $G_L$  be the subgraph induced by  $V_L$ . Then we define  $T_L$  as the number of triangles in  $G_L$ , and  $T_H = T - T_L$ .

We will compute estimates  $\overline{T}_L, \overline{T}_H$  of  $T_L, T_H$ , and estimate  $T$  as  $\overline{T} = \overline{T}_L + \overline{T}_H$ .

#### A.3.3 Estimating $T_L$

We will estimate  $T_L$  by sampling the vertices of  $V_L$  with probability  $\frac{1}{\sqrt{k}}$  each and calculating the number of triangles in the resulting graph.



Let  $c : V \rightarrow \{0, 1\}$  be a random hash function such that  $\forall v \in V$ :

$$c(v) = \begin{cases} 1 & \text{With probability } \frac{1}{\sqrt{k}}. \\ 0 & \text{Otherwise.} \end{cases}$$

Let  $V'_L = \{v \in V_L | c(v) = 1\}$ , and let  $G'_L$  be the subgraph of  $G$  induced by  $V'_L$ . Then  $\overline{T}_L$  is the number of triangles in  $G'_L$ , multiplied by  $k^{\frac{3}{2}}$ .

### A.3.4 Estimating $T_H$

We will estimate  $T_H$  by considering every vertex in  $V \setminus V_L$  separately. We will achieve this by sampling vertices  $v$  with probability  $2^{-\mathcal{T}_v}$ , and then sampling edges incident to  $v$  with probability proportional to  $2^{\mathcal{T}_v}$ .

Let  $h : V \rightarrow \{0, 1\}$  be a random hash function such that  $\forall v \in V$ :

$$h(v) = \begin{cases} i & \text{With probability } \frac{1}{\omega 2^i} \text{ for each } i \in \{0, \dots, \lceil \log \frac{\sqrt{k}}{\omega} \rceil\}. \\ -\infty & \text{Otherwise.} \end{cases}$$

And let  $r_C : V \rightarrow [0, 1]$  be a uniformly random hash function.

Then, for each  $v \in V_H = V \setminus V_L$ , we allow  $v$  to contribute to  $\overline{T}_H$  iff  $h(v) = \mathcal{T}_v$ . If it does, we calculate its contribution  $\overline{T}_v$  as follows:

We count a triangle  $t = (u, v, w)$  iff:

$$\begin{aligned} c(u) &= 1 \\ c(w) &= 1 \\ r_C(vu) &< \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}} \\ r_C(vw) &< \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}} \end{aligned}$$

Then, for each such triangle, we add  $\frac{k^2}{\omega 2^{\mathcal{T}_v}} \times \frac{1}{|\{x \in \{u, v, w\} | \mathcal{T}_x \neq \text{LIGHT}\}|}$  to  $\overline{T}_v$ . (with the second term then compensating for the fact that a triangle could potentially be counted at multiple different vertices)

We then define  $\overline{T}_H = \sum_{v \in V_H, h(v) = \mathcal{T}_v} \overline{T}_v$ .

## A.4 Final Output

We then output our estimate  $\overline{T}$  of  $T$  as  $\overline{T} = \overline{T}_L + \overline{T}_H$ .

# B Analysis

## B.1 First Pass

The following lemmas will bound the deviation of  $\mathcal{T}_v$  from its desired value— $\log \frac{\tilde{T}}{\omega T_v}$  if  $T_v \geq \frac{\tilde{T}}{\sqrt{k}}$ , and LIGHT otherwise.

**Lemma 18.** For any  $v \in V(G)$ , and for  $i \in [2]$ , let  $X_{v,i}^{(h)}$  be as defined previously. Then,  $\mathbb{E} \left[ X_{v,i}^{(h)} \right] = T_v \frac{\omega^2 2^{2h}}{k^2}$  and  $\text{Var} \left( X_{v,i}^{(h)} \right) \leq T_v \frac{\omega^2 2^{2h}}{k^2} + \sum_w T_{vw}^2 \frac{\omega^3 2^{3h}}{k^3}$ .

*Proof.* For each triangle  $t = (u, v, w)$ ,  $X_{v,i}^{(h,t)} = 1$  iff  $d_i(u) = d_i(w) = 1$  and  $r_{D,i}(vu), r_{D,i}(vw) < \frac{\omega 2^h}{\sqrt{k}}$ , which happens with probability  $\frac{\omega^2 2^{2h}}{k^2}$ .  $X_{v,i}^{(h)} = \sum_{t \in \tau(v)} X_{v,i}^{(h,t)}$ , so  $\mathbb{E} \left[ X_{v,i}^{(h)} \right] = T_v \frac{\omega^2 2^{2h}}{k^2}$ .

Then, to bound the variance, we start by bounding  $\mathbb{E} \left[ \left( X_{v,i}^{(h)} \right)^2 \right] = \sum_{t_1, t_2 \in \tau(v)} \mathbb{E} \left[ X_{v,t_1}^{(h)} X_{v,t_2}^{(h)} \right]$ .

We split the terms  $\mathbb{E} \left[ X_{v,t_1}^{(h)} X_{v,t_2}^{(h)} \right]$  by the value of  $l = |V(t_1) \cap V(t_2)|$ . (noting that  $l \in \{1, 2, 3\}$ )

Then,  $X_{v,t_1}^{(h)} X_{v,t_2}^{(h)} = 1$  implies that, for the  $5 - l$  vertices  $u \in V(t_1) \cup V(t_2) \setminus \{v\}$ ,  $d(u) = 1$ , and for the  $5 - l$  edges  $e \in \{vu | u \in V(t_1) \cup V(t_2) \setminus \{v\}\}$ ,  $r_D(e) < \frac{\omega 2^h}{\sqrt{k}}$ . So  $\mathbb{E} \left[ X_{v,t_1}^{(h)} X_{v,t_2}^{(h)} \right] \leq \left( \frac{\omega 2^h}{k} \right)^{5-l}$ .

There are no more than  $T_v^2$  such pairs for  $l = 1$ , no more than  $\sum_w T_{vw}^2$  for  $l = 2$ , and exactly  $T_v$  for  $l = 3$ , which gives us:

$$\begin{aligned} \text{Var} \left( X_{v,i}^{(h)} \right) &= \mathbb{E} \left[ \left( X_{v,i}^{(h)} \right)^2 \right] - \mathbb{E} \left[ X_{v,i}^{(h)} \right]^2 \\ &\leq T_v^2 \frac{\omega^4 2^{4h}}{k^4} + \sum_w T_{vw}^2 \frac{\omega^3 2^{3h}}{k^3} + T_v \frac{\omega^2 2^{2h}}{k^2} - \left( T_v \frac{\omega^2 2^{2h}}{k^2} \right)^2 \\ &= T_v \frac{\omega^2 2^{2h}}{k^2} + \sum_w T_{vw}^2 \frac{\omega^3 2^{3h}}{k^3} \end{aligned}$$

□

**Corollary 19.**  $\mathbb{P} \left[ X_{v,i}^{(h)} \leq \frac{T_v}{2} \frac{\omega^2 2^{2h}}{k^2} \right] \lesssim \frac{k^2}{\omega^2 2^{2h} T_v} + \sum_w \frac{T_{vw}^2}{T_v^2} \frac{k}{\omega 2^h}$ .

*Proof.*  $X_{v,i}^{(h)} \leq \frac{T_v}{2} \frac{\omega^2 2^{2h}}{k^2}$  implies that  $\left| X_{v,i}^{(h)} - \mathbb{E} \left[ X_{v,i}^{(h)} \right] \right| \geq \frac{1}{2} \mathbb{E} \left[ X_{v,i}^{(h)} \right]$ , and so by Chebyshev's inequality it occurs with probability  $\lesssim \frac{\text{Var} \left( X_{v,i}^{(h)} \right)}{\mathbb{E} \left[ X_{v,i}^{(h)} \right]^2}$ . □

**Lemma 20.** If  $T_v \geq 2 \frac{\tilde{T}}{\sqrt{k}}$ , then  $\mathbb{P} [\mathcal{T}_v = \text{LIGHT}] \lesssim \frac{k}{T_v} + \sum_w T_{vw}^2 \frac{\sqrt{k}}{T_v^2}$ .

*Proof.*  $\mathcal{T}_v = \text{LIGHT}$  implies that  $X_{v,i}^{(h)} < \frac{\omega 2^h \tilde{T}}{k^2}$  for all  $h \in \left\{ 1, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil \right\}$  and  $i \in [2]$ , and so in particular  $X_{v,1}^{\left( \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil \right)} < \frac{\tilde{T}}{k^{\frac{3}{2}}} \leq \frac{T_v}{2k} \leq \frac{T_v}{2} \frac{\omega^2 2^{2 \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil}}{k^2}$ .

So the result follows by using Corollary 19 with  $h = \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil$ . □

**Lemma 21.**  $\forall v \in V, \mathbb{P} [\mathcal{T}_v \neq \text{LIGHT}] \lesssim \frac{\sqrt{k} T_v}{\tilde{T}}$ .

*Proof.*  $\mathcal{T}_v \neq \text{LIGHT}$  implies that there is at least one  $h \in \left\{ 1, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil \right\}$  and  $i \in [2]$  such that  $X_{v,i}^{(h)} \geq \frac{\omega 2^h \tilde{T}}{k^2}$ . By Lemma 18,  $\mathbb{E} \left[ X_{v,i}^{(h)} \right] = T_v \frac{\omega^2 2^{2h}}{k^2}$ , so by Markov's inequality this occurs with probability  $\leq \frac{T_v}{\tilde{T}} \omega 2^h$ . So the probability that it holds for any  $h, i$  is  $\leq 2 \sum_{h=1}^{\left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil} \frac{T_v}{\tilde{T}} \omega 2^h \lesssim \sum_{i=0}^{\infty} \frac{T_v \sqrt{k}}{\tilde{T}} 2^{-i} \lesssim \frac{\sqrt{k} T_v}{\tilde{T}}$ . □

**Lemma 22.**  $\forall v \in V, l \geq 2, \mathbb{P} \left[ \mathcal{T}_v > \left\lceil \log \frac{\tilde{T}}{\omega T_v} \right\rceil + l \right] \lesssim \left( \frac{k^2 T_v}{2^{2l} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^l \tilde{T} T_v} \right)^2$ .

*Proof.*  $\mathcal{T}_v > \left\lceil \log \frac{\tilde{T}}{\omega T_v} \right\rceil + l$  implies that  $\forall i \in [2], X_{v,i}^{\left(\left\lceil \log \frac{\tilde{T}}{\omega T_v} \right\rceil + l\right)} < \frac{2^l (\tilde{T})^2}{k^2 T_v}$ . By applying Corollary 19 with  $h = \left\lceil \log \frac{\tilde{T}}{\omega T_v} \right\rceil + l$ , and using the independence of  $X_{v,1}^{(h)}, X_{v,2}^{(h)}$ , this happens with probability  $\lesssim \left( \frac{k^2 T_v}{2^{2l} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^l \tilde{T} T_v} \right)^2$ .  $\square$

## B.2 Second Pass

Making use of the lemmas from the previous section, we bound the error of our estimators  $\overline{T}_L$  and  $\overline{T}_H$ , and thus of our final estimator  $\overline{T}$ .

### B.2.1 $\overline{T}_L$

**Lemma 23.**  $\mathbb{E} [\overline{T}_L | T_L] = T_L$ .

*Proof.*  $T_L$  is the number of triangles in the random graph  $G_L \subseteq G$ , and  $\overline{T}_L$  is  $k^{\frac{3}{2}}$  times the number of triangles in the random graph  $G'_L \subseteq G_L$ . Any triangle  $t = (u, v, w)$  in  $G_L$  is in  $G'_L$  iff  $c(u) = c(v) = c(w) = 1$ , which occurs with probability  $\frac{1}{k^{\frac{3}{2}}}$ . So the expected number of triangles in  $G'_L$  is  $\frac{T_L}{k^{\frac{3}{2}}}$ , and so  $\mathbb{E} [\overline{T}_L | T_L] = T_L$ .  $\square$

**Lemma 24.**  $\text{Var}(\overline{T}_L) \lesssim T k^{\frac{3}{2}} + k T_E + T \tilde{T}$ .

*Proof.* Let  $T'_L$  be the number of triangles in  $G'_L$  (so  $\overline{T}_L = k^{\frac{3}{2}} T'_L$ ). A triangle  $t = (u, v, w)$  in  $G$  is in  $G'_L$  iff  $c(u) = c(v) = c(w)$  and  $\mathcal{T}_u = \mathcal{T}_v = \mathcal{T}_w = \text{LIGHT}$ . We proceed by bounding  $\mathbb{E} [T'^2_L]$ . This will be equal to the sum over every ordered pair  $(t_1, t_2)$  of the probability that both  $t_1$  and  $t_2$  are in  $G'_L$ . There are four scenarios to consider:

$t_1 = t_2 = (u, v, w)$ : Then both will be counted iff  $t_1$  is, which requires  $c(u) = c(v) = c(w)$ , so this happens with probability  $\leq \frac{1}{k^{\frac{3}{2}}}$ . There are no more than  $T$  such “pairs” in  $G_L$ , and so they contribute at most  $\frac{T}{k^{\frac{3}{2}}}$  to  $\mathbb{E} [T'^2_L]$ .

$t_1$  and  $t_2$  share one edge  $vw$ : Let  $u_1, u_2$  be the remaining two vertices. Then for both  $t_1$  and  $t_2$  to be in  $G'_L$ , it is necessary (but not sufficient) that:

$$\begin{aligned} c(u_1) &= 1 \\ c(u_2) &= 1 \\ c(v) &= 1 \\ c(w) &= 1 \end{aligned}$$

So this holds with probability at most  $\frac{1}{k^2}$ . There are no more than  $T_e^2$  such pairs in for each edge  $e$  in  $G_L$ , so they contribute at most  $\sum_e \frac{T_e^2}{k^2} = \frac{T_E}{k^2}$  to  $\mathbb{E} [T'^2_L]$  in total.

$t_1$  and  $t_2$  share only  $v$ : Let  $u_1, w_1, u_2, w_2$  be the remaining four vertices. Then for both  $t_1$  and  $t_2$  to be in  $G'_L$ , it is necessary that:

$$\begin{aligned} c(v) &= 1 \\ c(u_1) &= 1 \\ c(w_1) &= 1 \\ c(u_2) &= 1 \\ c(w_2) &= 1 \end{aligned}$$

So this holds with probability at most  $\frac{1}{k^{\frac{3}{2}}}$ . We now split such pairs into two categories.

$T_v \leq \frac{2\tilde{T}}{\sqrt{k}}$ : There are at most  $\sum_{v, T_v \leq \frac{2\tilde{T}}{\sqrt{k}}} T_v^2 \lesssim \frac{T\tilde{T}}{\sqrt{k}}$  such pairs in  $G_L$ , and so this category contributes at most  $\frac{T\tilde{T}}{k^3}$  to  $\mathbb{E}[T_L'^2]$ .

$T_v > \frac{2\tilde{T}}{\sqrt{k}}$ : In this case, we also use the bound on the probability that  $\mathcal{T}_v = \text{LIGHT}$  from Lemma 20. This is  $\lesssim \frac{k}{T_v} + \sum_w T_{vw}^2 \frac{\sqrt{k}}{T_v^2}$ , so the contribution to  $\mathbb{E}[T_L'^2]$  from this category is  $\lesssim \sum_v \frac{T_v^2}{k^{\frac{3}{2}}} \left( \frac{k}{T_v} + \sum_w T_{vw}^2 \frac{\sqrt{k}}{T_v^2} \right) = \frac{T}{k^{\frac{3}{2}}} + \frac{T_E}{k^2}$ .

$t_1$  and  $t_2$  are disjoint: Then for  $t_1$  and  $t_2$  to be in  $G'_L$ , it is necessary that all 6 vertices  $v$  of  $t_1, t_2$  have  $c(v) = 1$ , which happens with probability  $\frac{1}{k^3}$ . So then, letting  $p_t$  be the probability that a triangle  $t \in G$  is in  $G_L$ , the total contribution to  $\mathbb{E}[T_L'^2]$  from this case will be at most  $\sum_{t_1 \cap t_2 = \emptyset} \frac{p_{t_1} p_{t_2}}{k^3} \leq \frac{1}{k^3} \mathbb{E}[T_L]^2 = \mathbb{E}[T_L']^2$ .

By summing these together, as  $\overline{T}_L = k^{\frac{3}{2}} T_L'$ , this gives us:

$$\begin{aligned} \text{Var}(\overline{T}_L) &= k^3 \text{Var}(T_L') \\ &= k^3 (\mathbb{E}[T_L'^2] - \mathbb{E}[T_L']^2) \\ &\lesssim k^3 \left( \frac{T}{k^{\frac{3}{2}}} + \frac{T_E}{k^2} + \frac{T\tilde{T}}{k^3} + \mathbb{E}[T_L']^2 - \mathbb{E}[T_L']^2 \right) \\ &\lesssim T k^{\frac{3}{2}} + k T_E + T\tilde{T} \end{aligned}$$

□

### B.2.2 $\overline{T}_H$

**Lemma 25.**  $\mathbb{E}[\overline{T}_H | T_H] = T_H$ .

*Proof.*  $\overline{T}_H = \sum_{v \in V_H} \overline{T}_v$ , and  $T_H$  is the number of triangles in  $G$  that are not in  $G_L$ . If a triangle in  $G$  is in  $G_L$ , it can never contribute to  $\overline{T}_v$  for any  $v$ , as every vertex  $v$  it uses will have  $\mathcal{T}_v = \text{LIGHT}$ . If a triangle in  $G$  is not in  $G_L$ , it has  $s \in \{1, 2, 3\}$  vertices  $u$  s.t.  $\mathcal{T}_u \neq \text{LIGHT}$ . It will therefore have  $s$  vertices  $v$  such that it can contribute to  $T_v$ .

At each of those vertices  $v$ , if the triangle is  $t = (u, v, w)$ , it will contribute  $\frac{k^2}{\omega^{2\mathcal{T}_v s}}$  to  $\overline{T}_v$  iff:

$$\begin{aligned}
h(v) &= \mathcal{T}_v \\
c(u) &= 1 \\
c(w) &= 1 \\
r_C(vu) &< \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}} \\
r_C(vw) &< \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}}
\end{aligned}$$

This happens with probability  $\frac{1}{\omega 2^{\mathcal{T}_v}} \times \frac{1}{\sqrt{k}} \times \frac{1}{\sqrt{k}} \times \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}} \times \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}} = \frac{\omega 2^{\mathcal{T}_v}}{k^2}$ , so the expected contribution of  $t$  to  $\overline{T}_v$  is  $\frac{1}{s}$ . Therefore, as there are  $s$  vertices where  $t$  can contribute, its expected contribution to  $\overline{T}_H$  is 1.

Therefore,  $\mathbb{E}[\overline{T}_H | T_H]$  is precisely the number of triangles in  $G$  that are not in  $G_L$ , which is  $T_H$ .  $\square$

**Lemma 26.**  $\text{Var}(\overline{T}_H) \lesssim T \frac{k^2}{\omega} + T_E k + T^2$ .

*Proof.* We now care, for any triangle  $t$ , which vertex we are counting it at (and so which  $T_v$  it may contribute to). We will therefore use  $t^v$  to denote the triangle  $t$ , counted at the vertex  $v$ . We then define  $Y_{t^v}$  as  $\frac{k^2}{\omega 2^{\mathcal{T}_v}}$  if  $t$  is counted at  $v$  and 0 otherwise.

Then,  $\overline{T}_v \leq \sum_{t \in \tau(v)} Y_{t^v}$ , as  $Y_{t^v}$  is equal or greater to the contribution to  $T_v$  from  $t$ . So with  $Y = \sum_v \sum_{t \in \tau(v)} Y_{t^v}$ ,  $\mathbb{E}[\overline{T}_H^2] \leq \mathbb{E}[Y^2]$ . We will bound  $\mathbb{E}[Y^2]$  by bounding  $\mathbb{E}[Y_{t_1^u} Y_{t_2^v}]$  for each pair  $t_1^u, t_2^v$ .

We will bound the total contribution to  $Y^2$  from each possible value of  $l = |V(t_1) \cap V(t_2)|$ , treating  $l = 1$  as a special case, and treating  $u = v$  and  $u \neq v$  separately.

$l \neq 1, u \neq v$ :  $Y_{t^v} > 0$  implies that  $h(u) = \mathcal{T}_u, h(v) = \mathcal{T}_v$ , which happens with probability  $\omega^{-2} 2^{-\mathcal{T}_u - \mathcal{T}_v}$ .

It also implies that  $\forall w \in (V(t_1) \setminus \{u\}) \cup (V(t_2) \setminus \{v\}), c(w) = 1$ . There are  $4 - l$  vertices in this set, so this happens with probability  $\leq k^{-2 + \frac{l}{2}}$ . It also implies that  $\forall e \in \{uw | uw \in E(t_1)\}, r_C(e) < \frac{\omega 2^{\mathcal{T}_u}}{\sqrt{k}}$  and  $\forall e \in \{vw | vw \in E(t_2)\}, r_C(e) < \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}}$ . These sets can overlap in at most one edge  $(uv)$ , and each have size 2, so this happens with probability  $\leq \left(\frac{\omega}{\sqrt{k}}\right)^3 2^{2\mathcal{T}_u + 2\mathcal{T}_v - \max\{\mathcal{T}_u, \mathcal{T}_v\}}$ . Furthermore, the overlap in  $uv$  can only occur if  $u \in V(t_2)$  and  $v \in V(t_1)$ , and in this case we will also need  $c(u) = 1$  and  $c(v) = 1$ . So this either reduces the probability by a factor of  $\frac{\omega 2^{\max\{\mathcal{T}_u, \mathcal{T}_v\}}}{\sqrt{k}}$  or  $\frac{1}{k}$ , and so in either case by at least a factor of  $\frac{\omega 2^{\max\{\mathcal{T}_u, \mathcal{T}_v\}}}{\sqrt{k}}$ .

So as these three conditions are independent, and multiplying the probability that  $Y_{t_1^u}, Y_{t_2^v} > 0$  by the values they take when they are,  $\mathbb{E}[Y_{t_1^u} Y_{t_2^v}] \leq k^{\frac{l}{2}}$ . So the contribution to  $\mathbb{E}[Y^2]$  from such pairs is  $\lesssim T^2$  for  $l = 0$  (as there are at most  $T^2$  such pairs in  $G$ ),  $T_E k$  for  $l = 2$  (as there are at most  $\sum_e T_e^2 = T_E$  such pairs in  $G$ ), and  $T k^{\frac{3}{2}}$  for  $l = 3$  (as any such ‘‘pair’’ has  $t_1 = t_2$ , so there are at most  $T$  of them).

$l \neq 1, u = v$ : Note that as  $u = v, l \neq 0$ .  $Y_{t^v} > 0$  implies that  $h(v) = \mathcal{T}_v$ , which happens with probability  $\omega^{-1} 2^{-\mathcal{T}_v}$ . It also implies that  $\forall w \in V(t_1) \setminus \{u\} \cup V(t_2) \setminus \{v\}, c(w) = 1$ . There are  $5 - l$  vertices in this set, so this happens with probability  $\leq k^{-\frac{5-l}{2}}$ . It also implies that

$\forall e \in \{uw|uw \in E(t_1)\} \cup \{vw|vw \in E(t_2)\}$ ,  $r_C(e) < \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}}$ . As this set has  $5 - l$  elements, this happens with probability  $\leq \left(\frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}}\right)^{5-l}$ .

So as these three conditions are independent, and multiplying the probability that  $Y_{t_1^u}, Y_{t_2^v} > 0$  by the values they take when they are,  $\mathbb{E}[Y_{t_1^u} Y_{t_2^v}] \leq (\omega 2^{\mathcal{T}_v})^{2-l} k^{l-1}$ . So as  $\omega 2^{\mathcal{T}_v} \leq \sqrt{k}$ , this is  $\leq k^{\frac{l}{2}}$  for  $l < 3$ , and  $\leq k^{l-1}$  otherwise. So the contribution to  $\mathbb{E}[Y^2]$  from such pairs is  $T_E k$  from  $l = 2$  and  $T k^{\frac{3}{2}}$  from  $l = 3$ .

$l = 1, u \neq v$ : As in the  $l \neq 1$  case,  $\mathbb{E}[Y_{t_1^u} Y_{t_2^v}] \leq k^{\frac{l}{2}} = \sqrt{k}$ . So we seek to bound the number of such pairs.

For any  $w \in V(G)$ , the number of such pairs intersecting at  $w$  is  $\leq \left(\sum_{v \in V(G), \mathcal{T}_v \neq \text{LIGHT}} T_{wv}\right)^2$ . Now let  $L = |\{v \in V(G) | \mathcal{T}_v \neq \text{LIGHT}\}|$ . Suppose  $L = r$ . By Cauchy-Schwartz,

$$\left(\sum_{v \in V(G), \mathcal{T}_v \neq \text{LIGHT}} T_{wv}\right)^2 \leq r \sum_{v \in V(G)} T_{wv}^2$$

. So by summing across all  $w$ , the total number of such pairs is  $\leq r T_E$ . So the contribution to the expectation conditioned on  $L = r$  is  $\leq \sqrt{k} r T_E$ . As our bound on  $\mathbb{E}[Y_{t_1^u} Y_{t_2^v}]$  holds for any values of  $\mathcal{T}_u, \mathcal{T}_v$ , we can then bound the contribution to  $\mathbb{E}[Y^2]$  from such pairs by:

$$\begin{aligned} \sum_r \sqrt{k} r T_E \mathbb{P}[L = r] &= \mathbb{E}[L] \sqrt{k} T_E \\ &= \sqrt{k} T_E \sum_{v \in V(G)} \mathbb{P}[\mathcal{T}_v \neq \text{LIGHT}] \\ &\lesssim \sqrt{k} T_E \sum_{v \in V(G)} \frac{\sqrt{k} T_v}{\tilde{T}} && \text{By Lemma 21.} \\ &= T_E k \frac{T}{\tilde{T}} \end{aligned}$$

$l = 1, u = v$ : As in the  $l \neq 1$  case,  $\mathbb{E}[Y_{t_1^u} Y_{t_2^v}] \leq \frac{k^{l-1}}{\omega^{l-2} 2^{\mathcal{T}_v(l-2)}} = \omega 2^{\mathcal{T}_v}$ . So at any vertex  $v$ , the contribution to the expectation, conditioned on  $\mathcal{T}_v$ , is  $\leq T_v^2 \omega 2^{\mathcal{T}_v}$ .

We will consider two cases:  $\mathcal{T}_v \leq \left\lceil \log \frac{\tilde{T}}{\omega T_v} \right\rceil$  and  $\mathcal{T}_v = \left\lceil \log \frac{\tilde{T}}{\omega T_v} \right\rceil + i$  for some  $i \geq 1$ .

In the first case,  $\mathbb{E}[Y_{t_1^u} Y_{t_2^v}] \leq \frac{\tilde{T}}{T_v}$ , and there are no more than  $T_v^2$  such pairs for each vertex  $v$ , so the total contribution to the expectation from such vertices is  $\leq \sum_{v \in V(G)} T_v \tilde{T} \lesssim T \tilde{T}$ .

In the second case,  $\mathbb{E}[Y_{t_1^u} Y_{t_2^v}] \leq T_v \tilde{T} 2^i$ , and the probability of  $\mathcal{T}_v$  being at least this high is:

$$\lesssim \min \left\{ 1, \left( \frac{k^2 T_v}{2^{2i} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^i \tilde{T} T_v} \right)^2 \right\} \quad \text{By Lemma 22}$$

Now let  $x \in \mathbb{R}$  be the unique solution to  $\frac{k^2 T_v}{2^{2x} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^x \tilde{T} T_v} = 1$ . For  $i \leq x$ ,

$$1 \leq 2^{-|i-x|} \left( \frac{k^2 T_v}{2^{2i} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^i \tilde{T} T_v} \right)$$

. For  $i \geq x$ ,

$$\left( \frac{k^2 T_v}{2^{2i} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^i \tilde{T} T_v} \right)^2 \leq 2^{-|i-x|} \left( \frac{k^2 T_v}{2^{2i} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^i \tilde{T} T_v} \right)$$

So in either case, the probability of  $\mathcal{T}_v$  being at least this high is

$$\lesssim 2^{-|i-x|} \left( \frac{k^2 T_v}{2^{2i} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^i \tilde{T} T_v} \right)$$

. By summing across all possible values of  $i$ , it follows that the contribution to  $\mathbb{E}[Y^2]$  is:

$$\begin{aligned} &\lesssim T_v \tilde{T} \sum_{i=\max\{0, -\lceil \log \frac{\tilde{T}}{\omega T_v} \rceil\}}^{\lceil \log \frac{\sqrt{k}}{\omega} \rceil} 2^{i-|i-x|} \left( \frac{k^2 T_v}{2^{2i} (\tilde{T})^2} + \sum_w T_{vw}^2 \frac{k}{2^i \tilde{T} T_v} \right) \\ &\leq \sum_{i=\max\{0, -\lceil \log \frac{\tilde{T}}{\omega T_v} \rceil\}}^{\lceil \log \frac{\sqrt{k}}{\omega} \rceil} 2^{-|i-x|} \left( \frac{k^2 T_v^2}{\tilde{T}} + \sum_w T_{vw}^2 k \right) \\ &\leq \sum_{i=0}^{\infty} 2^{-i} \left( \frac{k^2 T_v^2}{\tilde{T}} + \sum_w T_{vw}^2 k \right) \\ &\lesssim \frac{k^2 T_v^2}{\tilde{T}} + \sum_w T_{vw}^2 k \end{aligned}$$

And so, by summing across all  $v$ , the contribution is:

$$\lesssim \frac{k^2 T_V}{\tilde{T}} + T_E k$$

Now, by summing the bounds from all of the above cases together, we can bound  $\mathbb{E}[Y^2]$  and therefore  $\text{Var}(\overline{T_H})$ .

$$\begin{aligned} \text{Var}(\overline{T_H}) &\leq \mathbb{E}[\overline{T_H}^2] \\ &\leq \mathbb{E}[Y^2] \\ &\lesssim T^2 + T\tilde{T} + \frac{k^2 T_V}{\tilde{T}} + T_E k + T k^{\frac{3}{2}} \end{aligned}$$

□

### B.2.3 $\bar{T}$

**Lemma 27.**  $\mathbb{E}[\bar{T}] = T$ .

*Proof.* For any  $T_L, T_H$ ,  $T_L + T_H = T$ . So as  $\bar{T} = \bar{T}_L + \bar{T}_H$ , by Lemmas 23, 25,  $\mathbb{E}[\bar{T}|T_L, T_H] = T_L + T_H = T$ .  $\square$

**Lemma 28.**  $\text{Var}(\bar{T}) \lesssim Tk^{\frac{3}{2}} + \frac{k^2 T_V}{T} + T_E k + T^2 + T\tilde{T}$ .

*Proof.*  $\bar{T} = \bar{T}_L + \bar{T}_H$ , so  $\text{Var}(\bar{T}) \lesssim \text{Var}(\bar{T}_L) + \text{Var}(\bar{T}_H)$ . Our bound then follows by using the bounds in Lemmas 24 and 26.  $\square$

## C Single-Pass Algorithm

### C.1 Outline

The algorithm presented earlier calculates an estimate of  $T$  in two conceptual “passes”. We show how, with only one pass, we can calculate the output of this algorithm while storing only

$O\left(\frac{m \log \frac{T_V \sqrt{k}}{T^2}}{k}\right)$  edges. Our main theorem then follows as a corollary.

We can do this because, in order to calculate the output of the second pass, we only need to know  $\mathcal{T}_v$  when it is equal to  $h(v)$ . In order to calculate  $\mathcal{T}_v$ , we need all the edges  $uv$  such that  $r_{D,i}(uv) < \frac{2^{h(u)}\omega}{\sqrt{k}}$  (for  $i = 1, 2$ ) and  $d(v) = 1$ . So as there are  $\sim \frac{\log \frac{T_V \sqrt{k}}{T^2}}{k}$  possible values  $h$  of  $h(u)$ , each of which is taken with probability  $\omega 2^{-h}$ , this means we need to store  $O\left(\frac{m \log \frac{T_V \sqrt{k}}{T^2}}{k}\right)$  edges in expectation.

We also show that the post-processing time required is  $O\left(\bar{m}^{\frac{3}{2}}\right)$ , where  $\bar{m}$  is the number of edges we sample.

### C.2 Space Complexity

**Lemma 29.** *The first and second pass calculations can be performed while storing no more than*

$O\left(\frac{m \log \frac{T_V \sqrt{k}}{T^2}}{k}\right)$  *edges in expectation.*

*Proof.* We define 5 sets of edges  $(E_i)_{i=0}^4$  that our algorithm will store. These are expressed as ordered pairs  $uv$ , but as our input is undirected edges, an edge  $e = uv$  will be kept if either  $uv$  or  $vu$  is in one of the sets  $E_i$ .

$$\begin{aligned} E_0 &= \{uv \in E \mid \exists i, d_i(u) = 1, d_i(v) = 1\} \\ E_1 &= \{uv \in E \mid \exists i, c(u) = 1, d_i(v) = 1\} \\ E_2 &= \{uv \in E \mid c(u) = 1, c(v) = 1\} \\ E_3 &= \{uv \in E \mid \exists i, r_{D,i}(uv) < \frac{2^{h(u)}\omega}{\sqrt{k}}, d_i(v) = 1\} \\ E_4 &= \{uv \in E \mid r_C(uv) < \frac{2^{h(u)}\omega}{\sqrt{k}}, c(v) = 1\} \end{aligned}$$



Then,  $\forall uv \in E$ :

$$\begin{aligned}
\mathbb{P}[uv \in E_0] &\lesssim \frac{1}{k} \\
\mathbb{P}[uv \in E_1] &\lesssim \frac{1}{k} \\
\mathbb{P}[uv \in E_2] &= \frac{1}{k} \\
\mathbb{P}[uv \in E_3] &= \frac{1}{\sqrt{k}} \sum_{i=1}^2 \sum_{h=0}^{\lceil \log \frac{\sqrt{k}}{\omega} \rceil} \mathbb{P} \left[ r_{D,i}(uv) < \frac{2^{h(u)} \omega}{\sqrt{k}} \mid h(u) = h \right] \mathbb{P}[h(u) = h] \\
&\lesssim \frac{1}{\sqrt{k}} \sum_{h=0}^{\lceil \log \frac{\sqrt{k}}{\omega} \rceil} \frac{2^h \omega}{\sqrt{k}} \frac{1}{\omega 2^h} \\
&\lesssim \frac{\log \frac{\sqrt{k}}{\omega}}{k} \\
\mathbb{P}[uv \in E_4] &\lesssim \frac{\log \frac{\sqrt{k}}{\omega}}{k}
\end{aligned}$$

So these sets will contain  $O\left(\frac{m \log \frac{\sqrt{k}}{\omega}}{k}\right)$  edges in expectation. Then, as  $\omega = \min\left\{\frac{(\tilde{T})^2}{T_V^+}, \sqrt{k}\right\}$ , and  $T_V^+$  can be any upper bound on  $T_V$ , maintaining these sets requires storing  $O\left(\frac{m \log \frac{T_V \sqrt{k}}{T^2}}{k}\right)$  edges.

We will then use these to calculate  $\overline{T}_L$  and  $\overline{T}_H$  as follows:

$\overline{T}_L$ : Recall that  $\overline{T}_L$  is the number of triangles in  $G'_L$ , the subgraph of  $G$  induced by  $V'_L = \{u \in V(G) \mid c(u) = 1, \mathcal{T}_u = \text{LIGHT}\}$ .  $E_2$  will contain every edge in  $E(G'_L)$ , so we can calculate  $\overline{T}_L$  provided we can determine which of the edges in  $E_2$  are in  $G'_L$ . It is therefore sufficient to calculate  $\mathcal{T}_v$  for all  $v$  s.t.  $c(v) = 1$ .

To do this, we will need to calculate  $\sum_{t \in \tau(v)} X_{v,i}^{(h,t)}$  for  $h = \lceil \log \frac{\sqrt{k}}{\omega} \rceil$ , and each  $i \in [2]$ . For each triangle  $(u, v, w) \in \tau(v)$ ,  $X_{v,i}^{(h,t)} = 1$  iff

$$\begin{aligned}
d_i(u) &= 1 \\
d_i(w) &= 1 \\
r_{D,i}(vu) &< \frac{\omega 2^h}{\sqrt{k}} \\
r_{D,i}(vw) &< \frac{\omega 2^h}{\sqrt{k}}
\end{aligned}$$

And 0 otherwise. The third and fourth conditions always hold, as we are dealing with the case when  $h = \lceil \log \frac{\sqrt{k}}{\omega} \rceil$ , so we need to know how many such triangles exist with  $d_i(u) = d_i(w) = 1$ .

When  $d_i(u) = d_i(w) = 1$ ,  $uw \in E_0$  and  $vu, vw \in E_1$  (as  $c(v) = 1$ ), so  $\sum_{t \in \tau(v)} X_{v,i}^{(h,t)}$  will be equal to the number of triangles  $(u, v, w)$  with  $d_i(u) = d_i(w) = 1$  in the edges we have sampled. So then  $\mathcal{T}_{v,i} = \text{LIGHT}$  iff this number is  $< \frac{\bar{T}}{k^{\frac{3}{2}}}$ .

So we can compute  $\mathcal{T}_v$  for each  $v$  s.t.  $c(v) = 1$  using our sampled edges, and therefore we can compute  $\overline{T}_L$ .

$\overline{T}_H$ : For any  $v \in V(G)$ ,  $\overline{T}_v = 0$  if  $h(v) \neq \mathcal{T}_v$ . So it is sufficient to calculate  $\mathcal{T}_v$  when  $h(v) = \mathcal{T}_v$ , and to know that  $h(v) \neq \mathcal{T}_v$  otherwise. We can calculate  $\sum_{t \in \tau(v)} X_{v,i}^{(h,t)}$  for each  $h \leq h(v)$ ,  $i \in [2]$ , as for any any triangle  $(u, v, w) \in G$  s.t.  $d_i(u) = d_i(w) = 1$ ,  $uw$  will be in  $E_0$ , and if  $r_{D,i}(vu), r_{D,i}(vw) < \frac{\omega 2^h}{\sqrt{k}}$  for some  $h \leq h(v)$ ,  $vu$  and  $vw$  will be in  $E_3$ .

So then, as  $H_{v,i} = \left\{ h \in \left\{ 1, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil \right\} \mid X_{v,i}^{(h)} \geq \frac{\omega \tilde{T} 2^h}{k^2} \right\}$  we can compute  $H_{v,i} \cap \{0, \dots, h(v)\}$ . As  $\mathcal{T}_{v,i} = \min H_{v,i}$ , we can compute each  $\mathcal{T}_{v,i}$  if it is  $\leq h$ , and determine that it is  $> h(v)$  or LIGHT otherwise. Then, if  $\mathcal{T}_v \leq h(v)$ , at least one of  $\mathcal{T}_{v,1}, \mathcal{T}_{v,2}$  is  $\leq h(v)$ , and so we can calculate it and therefore calculate  $\mathcal{T}_v$ , and if not we can determine that  $\mathcal{T}_v$  is either  $> h(v)$  or LIGHT. (although not necessarily which one)

So if  $\mathcal{T}_v \neq h(v)$ , we know  $\overline{T}_v = 0$ , and then if  $\mathcal{T}_v = h(v)$ , we need to know how many triangles  $(u, v, w)$  in  $G$  there are such that:

$$\begin{aligned} c(u) &= 1 \\ c(w) &= 1 \\ r_C(vu) &< \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}} \\ r_C(vw) &< \frac{\omega 2^{\mathcal{T}_v}}{\sqrt{k}} \end{aligned}$$

If these criteria hold, then  $uw \in E_2$ , and  $vu, vw \in E_4$ . So we can calculate  $\overline{T}_v$  by counting the number of triangles amongst our sampled edges such that these criteria hold.

We then calculate  $\overline{T}_H$  by summing  $\overline{T}_v$  for each  $v \in V(G)$ .

□

**Theorem 1** (Triangle estimation). *If  $\tilde{T} \leq T$ , our algorithm obtains an  $(\epsilon, \delta)$  approximation to  $T$  while keeping  $O\left(\frac{m \log \frac{1}{\delta}}{\epsilon^2} \left(\frac{1}{\tilde{T}^{2/3}} + \frac{\sqrt{\Delta_V} \log \tilde{T}}{\tilde{T}} + \frac{\Delta_E \log \tilde{T}}{\tilde{T}}\right)\right)$  edges. If  $\tilde{T} > T$ , the algorithm either determines that  $\tilde{T} > T$  or obtains a  $(1 \pm \epsilon)$  approximation to  $T$  with probability  $1 - \delta$ .*

*Proof.* We will assume that  $\epsilon \leq 1/2$ , as if  $\epsilon > 1/2$ , we may follow the procedure for an  $(1/2, \delta)$  approximation while losing at most a constant factor in the number of edges stored.

By Lemma 27, the expectation of our estimate  $\overline{T}$  of  $T$  will be  $T$ . We now choose  $k$  such that the variance of  $\overline{T}$  is  $O(T\tilde{T} + T^2)$ . By Lemma 28,  $\text{Var}(\overline{T}) \lesssim Tk^{\frac{3}{2}} + \frac{k^2 T_V}{T} + T_E k + T\tilde{T} + T^2$ . Then, as we may choose  $T_V^+$  to be any upper bound on  $T_V$ , we can choose it to be within a constant factor of the true value, and so the variance will be  $O(T\tilde{T} + T^2)$  if  $k = \min \left\{ \tilde{T}^{2/3}, \frac{\tilde{T}^{\frac{3}{2}}}{\sqrt{T_V}}, \frac{\tilde{T}^2}{T_E} \right\}$ .

We can repeat the algorithm  $O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta}\right)$  times. We split these results into blocks of size  $O\left(\frac{1}{\epsilon^2}\right)$ , take the mean of each block, and then take the median of all these means, which gives us an estimate

within  $\frac{1}{4}\epsilon \left( \sqrt{T\tilde{T}} + T \right)$  of  $T$  with probability  $1 - \delta$ , provided we choose the constant factors in our block sizes appropriately. We will then return this result if it is at least  $\frac{3}{4}\tilde{T}$ , and report that  $\tilde{T} > T$  otherwise.

To show that this gives a correct result (that is, giving an  $1 \pm \epsilon$  approximation if  $\tilde{T} < T$ , and either giving a  $1 \pm \epsilon$  approximation or reporting that  $\tilde{T} > T$  otherwise) with probability  $1 - \delta$ , we consider three cases.

$\tilde{T} \leq T$ : With probability  $1 - \delta$ , our result will be within  $\frac{\epsilon}{2}T$  of  $T$ . As  $\epsilon < 1/2$ , it will therefore be greater than  $\frac{3}{4}\tilde{T}$ , so we will return a  $1 \pm \epsilon/2$ , and therefore a  $1 \pm \epsilon$ , approximation of  $T$ .

$\tilde{T} \in (T, 2T)$ : With probability  $1 - \delta$ , our result will be within  $\epsilon T$  of  $T$ . So we will return a  $1 \pm \epsilon$  approximation of  $T$ , or report that  $\tilde{T} > T$ , both of which are valid results for  $\tilde{T}$  in this range.

$\tilde{T} > 2T$ : With probability  $1 - \delta$ , our result will be within  $\frac{\epsilon}{2}\tilde{T}$  of  $T$ . So as  $\epsilon < 1/2$ , it will be  $< \frac{3}{4}\tilde{T}$ , and so we will report that  $\tilde{T} > T$ .

By Lemma 29, we can obtain this approximation while keeping at most  $\lesssim \frac{m \log \frac{T_V \sqrt{k}}{\tilde{T}^2}}{k}$  edges. Then, if the dominating term in our expression for  $k$ ,  $\min \left\{ \tilde{T}^{2/3}, \frac{\tilde{T}^{\frac{3}{2}}}{\sqrt{T_V}}, \frac{\tilde{T}^2}{T_E} \right\}$ , is  $\tilde{T}^{2/3}$ ,  $\tilde{T}^{2/3} \leq \frac{\tilde{T}^{\frac{3}{2}}}{\sqrt{T_V}}$ , so  $\frac{T_V}{\tilde{T}^2} \leq \tilde{T}^{-\frac{1}{3}}$  and  $\sqrt{k} \leq \tilde{T}^{\frac{1}{3}}$ , so  $\frac{T_V \sqrt{k}}{\tilde{T}^2} \leq 1$ , and so we need  $\lesssim \frac{m}{k}$  edges. Otherwise, as  $\log \frac{T_V \sqrt{k}}{\tilde{T}^2} \leq \log k$ , we need  $\lesssim \frac{m \log k}{k}$  edges.

The result then follows from the fact that  $T_V \leq \Delta_V T$  and  $T_E \leq \Delta_E T$ .  $\square$

**Corollary 30.** *We can obtain a constant-error approximation to  $T$  while sampling  $\lesssim \frac{1}{\tilde{T}^{2/3}} + \frac{\sqrt{\Delta_V} \log \tilde{T}}{T} + \frac{\Delta_E \log \tilde{T}}{T}$  edges.*

### C.3 Time Complexity

**Lemma 31.** *If our algorithm samples  $\bar{m}$  edges, we can compute  $\bar{T}$  in  $O(\bar{m}^{\frac{3}{2}})$  time.*

*Proof.* Let  $(E_i)_{i=1}^5$  be as described in the previous section. After executing the sampling phase of our single-pass algorithm, by Lemma 29,  $E' = \bigcup_{i=1}^5 E_i$  contains every triangle we need to calculate  $\bar{T}$ . We can list the triangles in  $E'$  in  $O(\bar{m}^{\frac{3}{2}})$  time by using the algorithm in [IR78]. We can then calculate  $\bar{T}$  by passing over this list (of length no more than  $\bar{m}^{\frac{3}{2}}$ ) and, for each triangle in the list, checking whether it contributes to  $\bar{T}_L$  or  $\bar{T}_H$  and if so, what that contribution is.

Once the  $\mathcal{T}_v$  are known, this only requires checking a constant number of criteria per triangle. However, we must first calculate the  $\mathcal{T}_v$ . Doing this naively, by iterating over every triangle and checking all the  $X_{v,i}^h$  it could contribute to, could take as much as  $O(\bar{m}^{\frac{3}{2}} \log k)$  time.

Recall that a triangle  $(u, v, w)$  contributes to  $X_{v,i}^{(h)}$  iff  $d(u) = d(w) = 1$  and  $r_{D,i}(vu), r_{D,i}(vw) < \frac{2^h \omega}{\sqrt{k}}$ .  $X_{v,i}^{(h)}$  is therefore monotone increasing in  $h$ , so we can calculate all the  $X_{v,i}^{(h)}$  in  $O(\bar{m}^{\frac{3}{2}})$  time as follows:

1. Sort our edges  $e$  by the value of  $r_{D,i}(e)$ .
2. Copy each triangle  $(u, v, w)$  three times, choosing a different vertex as the “main vertex”  $v$  each time.

3. Sort these triangles in ascending order of  $\max\{r_{D,i}(vu), r_{D,i}(vw)\}$ . (and therefore by the least value of  $h$  needed for the triangle to contribute to  $X_{v,i}^{(h)}$ )
4. Run through the list in order, for each triangle adding 1 to  $X_{v,i}^{(h)}$  for each  $h$  such that

$$h \geq \log \left( \frac{\sqrt{k}}{\omega} \max\{r_{D,i}(vu), r_{D,i}(vw)\} \right)$$

By storing the running total as we run through the list, we can do this in a constant number of updates. For any  $v, i$ , the first time we set  $X_{v,i}^{(h)}$  to be greater than  $\frac{\omega \tilde{T}_2^h}{k^2}$  for some  $h$ , set  $\mathcal{T}_{v,i}$  to  $h$ .

For any  $v, i$  where  $\mathcal{T}_{v,i}$  has not been set at the end of this process, it is then LIGHT. We now have  $\mathcal{T}_{v,i}$  for  $i \in [2]$  and for every  $v$  involved in one of our triangles, so we may calculate  $\bar{T}$  in time  $O\left(\bar{m}i^{\frac{3}{2}}\right)$  by iterating over our list of triangles and checking the contribution of each one to  $\bar{T}_L$  and  $\bar{T}_H$ . □

## D Lower Bounds

We recall our definition of an *instance lower bound*.

**Definition 2.** *Let  $G$  be a graph. We say an algorithm solves  $G$  with  $S$  space/samples and  $(\epsilon, \delta)$  error if, for any  $G'$  isomorphic to some subgraph of  $G$ , the algorithm returns  $T(G') \pm \epsilon T(G)$  with  $1 - \delta$  probability, using no more than  $S$  space/samples.*

**Definition 3.** *For any given streaming model,  $\text{INSTOPT}(G, \epsilon, \delta)$  is the least amount of samples/space such that some algorithm solves  $G$  with  $\text{INSTOPT}(G, \epsilon, \delta)$  samples/space.*

For each bound, we will show two distributions on subgraphs that cannot be distinguished with less space/samples than our lower bound.

**Definition 6** (Distinguishing). *We say that an algorithm  $\mathcal{A}$  can distinguish two random graph distributions  $\mathcal{G}_1$  and  $\mathcal{G}_2$  if there exists  $f$  such that, for a pair of draws  $G_1$  and  $G_2$  from these distributions, and any relabelling of the vertices of  $G_1$  and  $G_2$ ,  $\Pr[f(\mathcal{A}(G_1)) = 1] \geq 3/4$  and  $\Pr[f(\mathcal{A}(G_2)) \neq 1] \geq 3/4$ .*

**Lemma 7.** *Let  $\mathcal{A}$  be an algorithm that solves triangle counting for a graph  $G$  with  $S$  space/samples and  $(\epsilon, 1/10)$  error. Then, for any two distributions  $\mathcal{G}_1, \mathcal{G}_2$  on subgraphs  $G_1$  and  $G_2$  of  $G$ , and  $C$  such that  $T(G_1) > C + \epsilon T(G)$  with  $\frac{9}{10}$  probability and  $T(G_2) < C - \epsilon T(G)$  with  $\frac{9}{10}$  probability,  $\mathcal{A}$  can distinguish them.*

*Proof.* By our definition of “solving” a graph, any algorithm which  $(\epsilon, 1/10)$  solves  $G$  can also distinguish between a draw from  $\mathcal{G}_1$  and  $\mathcal{G}_2$  by returning 1 for any graph with at least  $C$  triangles and 0 for any graph with fewer than  $C$  triangles, so that  $\Pr[f(\mathcal{A}(G_1)) = 1] \geq 4/5$  and  $\Pr[f(\mathcal{A}(G_2)) \neq 1] \geq 4/5$ . □

## D.1 Multiple Heavy Edges

We will express the bound shown in [BOV13] as an instance bound.

**Definition 8** (Heavy Edges Graph). *The heavy edges graph  $D_{r,d}$  consists of  $r$  copies of the following graph:  $d$  disjoint edges  $\{u_{2i}u_{2i+1}\}_{i=0}^{d-1}$ , one of which has both ends connected to a further  $d$  vertices  $\{v_i\}_{i=0}^{d-1}$ .*

**Theorem 9** (Heavy Edges Lower Bound). *For sufficiently small constant  $\epsilon$ ,  $\text{INSTOPT}(D_{r,d}, \epsilon, 1/10) = \Omega(d)$  bits in the insertion-only model.*

*Proof.* As in, [BOV13] we reduce to the indexing problem  $\text{Index}_n$ : Alice has a binary vector  $w$  of length  $n$ , and Bob has an index  $x \in [n]$ . Alice must send a message to Bob such that Bob can communicate  $w[x]$ . By [CCKM10], the randomized communication complexity of this problem is  $\Omega(n)$ . So for any instance of the problem:

**If  $r = 1$ :** Alice can encode  $w$  in a subgraph of  $D_{r,d}$  by, for each  $i$  in  $\{0, \dots, n-1\}$ , including  $u_{2i}u_{2i+1}$  iff  $w_i = 1$ . She can then run a distinguishing algorithm on these edges and send them to Bob, who proceeds to add the  $d$  vertices  $\{v_i\}_{i=0}^{d-1}$  to the graph, connecting each of them to  $u_{2x}$  and  $u_{2x+1}$ . If the algorithm reports the graph has 0 triangles, Bob reports that  $w_x = 0$ , and if it has  $d$  triangles, he reports that  $w_x = 1$ .

**If  $r > 1$ :** Alice and Bob can perform the same procedure but each copying their graph  $r$  times, with Bob now checking if he has  $rd$  triangles.

By [BOV13], there is a distribution on inputs such that this algorithm requires at least  $\Omega(d)$  bits of storage. Let  $G'_1, G'_2$  correspond to the distributions on the graph Alice and Bob create conditioned on  $w_x = 1$  and  $w_x = 0$  respectively. Note that  $T(G_1) = T(D_{r,d})$  and  $T(G_2) = 0$  with probability 1.

Then, letting  $H$  be an instance of  $D_{r,d}$ , each draw from a  $G'_i$  is isomorphic to a subgraph of  $H$ . We can therefore define subgraph distributions  $G_1, G_2$  on  $H$  and a permutation distribution  $\sigma$  such that  $\sigma^{-1}(G'_i) = G_i$ .

So we now have distributions on subgraphs  $G_1, G_2$  of  $H$  such that it is hard to distinguish between  $\sigma(G_1), \sigma(G_2)$  for some non-uniform permutation distribution  $\sigma$ . However, if it is hard to distinguish these, they must also be hard to distinguish for a uniformly random  $\sigma$ , as otherwise any distinguishing algorithm could distinguish them by randomly permuting its input.

So counting triangles for  $D_{r,d}$  requires  $\Omega(d)$  bits of storage. □

## D.2 Hubs Graph

**Definition 10** (Hubs Graph). *The hubs graph  $H_{r,d}$  consists of  $r$  copies of the following graph: a single vertex which participates in  $d$  edge-disjoint triangles.*

**Theorem 11** (Hubs Lower Bound). *For sufficiently small constant  $\epsilon$ ,  $\text{INSTOPT}(H_{r,d}, \epsilon, \delta) = \Omega(\sqrt{d})$  bits in the insertion-only model.*

*Proof.* We use a reduction from the variant of the Boolean Hidden Matching problem  $\text{BHM}_n$  set out in [KR06]. In this problem, Alice gets a string  $x \in \{0, 1\}^{2n}$ , Bob a perfect matching  $M$  on  $[2n]$  (interpreted as an  $n \times 2n$  matrix, where the  $n^{\text{th}}$  row corresponds to the  $n^{\text{th}}$  edge of the matching,

so that if the edge is  $ij$ , the  $i^{\text{th}}$  and  $j^{\text{th}}$  position of the row are 1, and all others are 0) and a string  $w \in \{0, 1\}^n$ .

Then, with  $h$  the hamming distance function, and matrix multiplication interpreted mod 2, Bob must determine whether  $h(Mx, w) \leq \frac{n}{3}$  or  $h(Mx, w) \geq \frac{2n}{3}$ .

In [KR06], it was shown that this requires  $\Omega(\sqrt{n})$  bits of communication in the randomized one-way communication model.

For each  $r \in \mathbb{N}$ , we provide a reduction from  $\text{BHM}_n$  to the problem of solving  $H_{r,n}$  as follows:

1. Alice encodes her string  $x \in \{0, 1\}^{2n}$  as a graph as follows: a single vertex  $a$ ,  $2n$  vertices  $\{b_i\}_{i \in [2n]}$ , and  $2n$  vertices  $\{c_i\}_{i \in [2n]}$ . Then, for each  $i \in [2n]$ , she adds the edge  $ab_i$  if  $x_i = 0$ , and  $uc_i$  if  $x_i = 1$ . She then creates  $r$  disjoint copies of this graph, runs the distinguishing algorithm on the whole graph, and sends the internal state of the algorithm to Bob.
2. Bob then, using the same vertex IDs, and the same number of copies, encodes  $M$  and  $w$  as follows: For the  $k^{\text{th}}$  edge  $ij$  in the matching, he adds  $b_i b_j, c_i c_j$  to his graph if  $w_k = 0$ , and  $b_i c_j, b_j c_i$  if  $w_k = 1$ . He then starts the triangle counting algorithm with the internal state sent to him by Alice, inserts the edges he has created (including all  $r$  disjoint copies), and reads the output of the algorithm. If the graph has  $\geq \frac{2rn}{3}$  triangles, he reports that  $h(Mx, w) \leq \frac{n}{3}$ , and if it has  $\leq \frac{rn}{3}$  triangles, he reports that  $h(Mx, w) \geq \frac{2n}{3}$ .

For each of the  $\frac{n}{T^2}$  copies Alice and Bob create, and for each  $r \in [n]$  (with the  $r^{\text{th}}$  edge of the matching being  $ij$ ), the graph will contain the triangle  $(a, b_i, b_j)$  if  $w_r = 0$  and  $x_i = x_j = 0$ ,  $(a, c_i, c_j)$  if  $w_r = 0$  and  $x_i = x_j = 1$ ,  $(a, b_i, c_j)$  if  $w_r = 1$  and  $x_i = 0, x_j = 1$ , and  $(a, b_j, c_i)$  if  $w_r = 1$  and  $x_i = 1, x_j = 0$ . So in each copy, there will be one triangle for the  $r^{\text{th}}$  edge  $ij$  iff  $w_r = x_i + x_j \pmod 2$ . So  $T = r(n - h(Mx, w))$ , and so this protocol correctly solves  $\text{BHM}_n$ .

By [KR06], there is a distribution on inputs such that this algorithm requires at least  $\Omega(\sqrt{d})$  bits of storage, and  $h(Mx, w)$  is guaranteed to be either less than  $\frac{n}{3}$  or at least  $\frac{2n}{3}$ . Let  $G'_1, G'_2$  correspond to the distributions on the graph Alice and Bob create conditioned on  $h(Mx, w) \leq \frac{n}{3}$  and  $h(Mx, w) \geq \frac{2n}{3}$  respectively. Note that  $T(G_1) \geq \frac{2rn}{3} > \frac{rn}{3} \geq T(G_2)$  with probability 1.

Then, each of the  $G'_i$  will always consist of a single hub vertex,  $2n$  other vertices,  $n$  edges from the hub to the other vertices, and a perfect matching on the other vertices. Therefore, letting  $H$  be an instance of  $H_{r,2n}$ , each draw from a  $G'_i$  is isomorphic to a subgraph of  $H$ . We can therefore define subgraph distributions  $G_1, G_2$  on  $H$  and a permutation distribution  $\sigma$  such that  $\sigma^{-1}(G'_i) = G_i$ .

So we now have distributions on subgraphs  $G_1, G_2$  of  $H$  such that it is hard to distinguish between  $\sigma(G_1), \sigma(G_2)$  for some non-uniform permutation distribution  $\sigma$ . However, if it is hard to distinguish these, they must also be hard to distinguish for a uniformly random  $\sigma$ , as otherwise any distinguishing algorithm could distinguish them by randomly permuting its input.

So counting triangles for  $H_{r,2n}$  requires  $\Omega(\sqrt{n})$  bits of storage. □

### D.3 Independent Triangles

**Definition 12** (Independent Triangles Graph). *The independent triangles graph  $I_n$  consists of  $n$  vertex-disjoint triangles.*

**Theorem 13** (Independent Triangles Lower Bound). *For sufficiently small constant  $\epsilon$ ,  $\text{INSTOPT}(I_n, \epsilon, 1/10) = \Omega\left(n^{\frac{1}{3}}\right)$  samples in the sampling model.*

If  $n$  is constantly bounded, the result holds automatically, so we will assume  $n \geq 80$ .

We start by proving a lemma that applies to all graphs.

**Lemma 32.** *Let  $G$  be a graph with  $m$  edges. Let  $\mathcal{A}$  be an algorithm which, when given any graph isomorphic to  $G$  as input, samples at least one triangle with constant probability.  $\mathcal{A}$  must sample  $\Omega\left(\frac{m}{T^{2/3}}\right)$  edges in expectation.*

*Proof.* As such algorithm must work for an arbitrary permutation of the vertices of  $G$ , we may assume without loss of generality that for any sampling algorithm we consider, for all potential triangles  $t$  in  $G$ , the probability  $p_t$  that  $t$  is sampled by the algorithm if it exists is the same, and likewise for all potential edges  $e$  in  $G$ . Then, we can bound  $p_t$  as follows: For any set of  $l$  vertices, let  $R$  be the number of edges amongst these  $l$  vertices that are sampled. Clearly  $\mathbb{E}[R] \leq l^2 p_e$ . Conditioned on  $R = r$ , it can be shown ([eh14]) that there are  $\lesssim r^{\frac{3}{2}}$  triangles amongst the edges sampled. So as there are  $\Omega(l^3)$  triangles amongst the  $l$  vertices,  $p_t$  conditioned on  $R = r$  is  $\lesssim \left(\frac{\sqrt{r}}{l}\right)^3$ . By Markov's inequality,  $R \lesssim l^2 p_e$  with constant probability (and we may choose any constant). Letting this event be  $A$ , we therefore have  $p_t \lesssim p_e^{\frac{3}{2}}$  conditioned on  $A$ . So  $p_e$  must be  $\Omega\left(\frac{1}{T^{2/3}}\right)$  if we want to sample at least one triangle with constant probability.  $\square$

Let  $I$  be an instance of  $I_n$ . As  $T(I) = n$ , and  $|E(I)| = 3n$ , we therefore need to sample  $\Omega(n^{\frac{1}{3}})$  edges if we are to sample at least one triangle with constant probability. We now need to demonstrate that a sampling algorithm that does not sample any triangles from  $I$  will also fail to distinguish between certain graphs with very different numbers of triangles. We will need the following lemma to proceed:

**Lemma 33.** *Let  $H$  be an acyclic graph, and let  $\chi$  be a random 2-coloring of  $V(H)$ . The following graphs are identically distributed:*

$$\begin{aligned} H_1 &= (V(H), \{uv \in E(H) \mid \chi(u) = \chi(v)\}) \\ H_2 &= (V(H), \{uv \in E(H) \mid \chi(u) \neq \chi(v)\}) \end{aligned}$$

*Proof.* As  $H$  is acyclic, it is a forest. Clearly the behavior of two distinct trees in the forest under any given coloring are independent of one another, so it suffices to prove that the result holds for  $H$  a tree. Choose a vertex  $v$  to act as the root of  $H$ . For all subtrees  $H'$  of  $H$ , with roots  $r'$ , let  $H'_1 = H' \cap H_1, H'_2 = H' \cap H_2$ . We will show that for  $i = 1, 2$ ,  $H'_i$  includes each edge in  $H'$  independently with probability  $\frac{1}{2}$ , and independently of value of  $\chi(r')$ .

We proceed by structural induction. Suppose the result holds for all proper subtrees of  $H'$ . Then let  $C$  be the set of children of  $r'$ . As the proper subtrees of  $H'$  do not share any vertices, the inclusion of any edge in  $H'_i$  is independent between any two distinct proper subtrees, and independent of the values  $\chi$  takes on  $C$ . The only remaining edges are the edges  $\{r'c \mid c \in C\}$ . Then, for any value of  $\chi(r')$ , the inclusion of each  $r'c$  happens with probability  $\frac{1}{2}$ , independently of each other and of the value of  $\chi(r')$ . So as the proper subtree edges are included independently of the value of  $\chi$  on  $C$ , they are included independently of the  $r'c$ . So every edge in  $E(H')$  is included with probability  $\frac{1}{2}$ , independently of one another and the value of  $\chi(r')$ .

So  $H_1, H_2$  are identically distributed.  $\square$

Now, with  $\chi$  a random 2-coloring of the vertices of  $I$ , we define two subgraphs  $I_1, I_2$  of  $I$  as follows:

$$\begin{aligned} I_1 &= (V(I), \{uv \in E(I) \mid \chi(u) = \chi(v)\}) \\ I_2 &= (V(I), \{uv \in E(I) \mid \chi(u) \neq \chi(v)\}) \end{aligned}$$

Note that  $I_1, I_2$  both have  $\frac{n}{2}$  edges in expectation.  $T(I_2)$  is always 0, and  $I_1$  will include every triangle independently with probability  $\frac{1}{4}$ . So as  $n \geq 80$ ,  $T(I_1) \geq \frac{n}{8} > 0 \geq T(I_2)$  with probability  $\geq \frac{9}{10}$ .

So for any sampling algorithm to sample a triangle from either  $I_1$  or  $I_2$  with constant probability, it must sample  $\Omega\left(m^{\frac{1}{3}}\right)$  edges in expectation. Suppose it does not. Then let  $\bar{I}$  be the set of edges from  $I$  it samples. (and so  $\bar{I}_i = I_i \cap \bar{I}$  is the set of edges of edges from  $I_i$  it samples) As  $I$  contains no cycle longer than a triangle,  $\bar{I}$  is acyclic. Therefore, by Lemma 33,  $\bar{I}_1, \bar{I}_2$  are identically distributed. So conditioned on failing to sample a triangle, the algorithm is incapable of distinguishing between  $I_1$  and  $I_2$ .

So to distinguish between  $I_1$  and  $I_2$  a sampling algorithm must sample  $\Omega\left(n^{\frac{1}{3}}\right)$  edges in expectation.

#### D.4 $G_{n,p}$

**Theorem 14** ( $G_{n,p}$  Lower Bound). *There exists a constant  $C$  such that, provided  $p \geq \frac{C}{n}$ , and for sufficiently small constant  $\epsilon$ ,  $\text{INSTOPT}(G_{n,p}, \epsilon, 1/10) = \Omega\left(\frac{1}{p}\right)$  samples.*

Let  $G$  be a random draw from  $G_{n,p}$ . We will use the same pair of colorings as in our proof of the  $\frac{m}{T^{2/3}}$  bound. However, as a graph drawn from  $G_{n,p}$  can contain cycles longer than triangles, we need some extra work to prove the subsampled graph is acyclic. We start by extending a result from [eh14].

**Lemma 34.** *Let  $G$  be a graph with  $m$  edges. There are at most  $(2m)^{\frac{1}{2}}$  length- $l$  cycles in  $G$ .*

*Proof.* As in [eh14], we define  $\forall v \in V(G), A(v) = \{u \in N(v), d(u) \geq d(v)\}$ .  $\forall w \in V(G), |A(w)| \leq \sqrt{m}$ , if  $|A(w)| > \sqrt{2m}$ , then there are  $> \sqrt{2m}$  vertices of degree  $> \sqrt{2m}$  in the graph, and thus  $> m$  edges. So  $|A(w)| \leq \sqrt{2m}$ .

We order each  $l$ -cycle as  $(v_0, \dots, v_{l-1})$  so that the  $v_{l-1}$  has the highest degree in the cycle. (there are two such orderings, we pick between them arbitrarily) So each cycle now corresponds to  $l$  directed edges  $(v_0v_1, v_1v_2, \dots, v_{l-2}v_{l-1})$ .

Suppose  $l$  is even. The cycle is then uniquely determined by the  $\frac{l}{2}$  edges  $\{v_{2i}v_{2i+1}\}_{i \in [\frac{l}{2}]}$  and their directions. So there are at most  $(2m)^{\frac{l}{2}}$   $l$ -cycles in  $G$ .

Now suppose  $l$  is odd. The cycle is then uniquely determined by the  $\frac{l-1}{2}$  edges  $\{v_{2i}v_{2i+1}\}_{i \in [\frac{l-1}{2}]}$  and their directions, plus the final vertex  $v_{l-1}$ . There are at most  $(2m)^{\frac{l-1}{2}}$  possible choices for the edges and directions. Furthermore, as  $v_{l-1}$  has the highest degree of any vertex in the cycle,  $v_{l-1} \in A(v_{l-2}) \cap A(v_0)$ . So as this set has size  $\leq \sqrt{2m}$ , there are at most  $\sqrt{2m}$  possibilities for  $v_{l-1}$ , and so there are at most  $(2m)^{\frac{l}{2}}$   $l$ -cycles in  $G$ .  $\square$

We can now show that a sampling algorithm that does not use  $\Omega\left(\frac{1}{p}\right)$  samples will, with probability bounded below by an arbitrary constant, fail to sample any cycle in the graph. Let  $\bar{G}$  be the subgraph of  $K_n$  that the algorithm would sample, and let  $q$  be the probability of sampling any given edge. (as in the independent triangles case, we will choose our graph permutation uniformly at random, so without loss of generality we can assume this is the same for all edges) Then  $\mathbb{E}[|E(\bar{G})|] \leq qn^2$ . So, by Lemma 34 and Markov's inequality, the number of  $l$ -cycles in  $\bar{G}$  is  $\leq q^{\frac{l}{2}}(2n)^l$  with constant probability. (for any arbitrarily small constant, so let us choose it to be  $\frac{1}{2}$ ) So, conditioned on this constant probability event  $A$ , the probability of any  $l$ -cycle being in



$\overline{G} \cap G$  is  $\leq p^l q^{\frac{1}{2}} (2n)^l$ . So, taking the union bound across all  $l$ , this probability is  $\leq \sum_{l=3}^{\infty} p^l q^{\frac{1}{2}} (2n)^l$ . So  $q$  must be  $\Omega\left(\frac{1}{n^2 p^2}\right)$  for at least one  $l$ -cycle to be found with probability  $\geq \frac{1}{2}$ .

Then, for  $\chi$  a random coloring of  $[n]$ , let  $G_1, G_2$  be defined as follows:

$$\begin{aligned} G_1 &= (V(G), \{uv \in E(G) \mid \chi(u) = \chi(v)\}) \\ G_2 &= (V(G), \{uv \in E(G) \mid \chi(u) \neq \chi(v)\}) \end{aligned}$$

Note that  $\mathbb{E}[T(G)] = \binom{n}{3} p^3$ .  $T(G_2)$  will always be 0. Then, there are at least  $\lceil \frac{n}{2} \rceil$  vertices with the same coloring under  $\chi$ . So we can take  $H \subseteq G_1$  such that all vertices of  $H$  have the same color and  $H$  is distributed as  $G_{\lceil \frac{n}{2} \rceil, p}$ . Then, let  $B$  be the number of triangles in  $H$ .  $\mathbb{E}[B] = \binom{\lceil n/2 \rceil}{3} p^3$ , and by [MV08],  $\frac{\text{Var}(B)}{\mathbb{E}[B]^2} = O\left(\frac{1}{n^3 p^3} + \frac{1}{n^2 p}\right)$ , so for  $p = \Omega\left(\frac{1}{n}\right)$ ,  $B$  will be a constant fraction of  $\mathbb{E}[T(G)] = O(n^3 p^3)$  with probability  $\frac{9}{10}$ .

Now, by Lemma 33, conditioned on a sampling algorithm failing to sample any cycles, the samples taken from  $G_1$  and  $G_2$  are identically distributed, in which case the algorithm can distinguish them with probability at most  $\frac{1}{2}$ . So  $q$  must be  $\Omega\left(\frac{1}{n^2 p^2}\right)$  to distinguish them with probability  $\geq \frac{3}{4}$ , and so as  $G_1, G_2$  each contain  $\Omega(n^2 p)$  edges in expectation, the algorithm must sample  $\Omega\left(\frac{1}{p}\right)$  edges in expectation.

## E Triangle-Dependent Sampling Lower Bound

**Definition 15.** Let  $\mathcal{A}$  be a sampling algorithm for counting triangles. We say that  $\mathcal{A}$  is a triangle-dependent sampling algorithm if, for all graphs  $G$ ,  $\mathcal{A}(G)$  depends only on the set of triangles sampled by  $\mathcal{A}$ .

**Definition 16.** For any graph  $G$ ,  $\epsilon > 0$ , let the vertices  $v \in V(G)$  be ordered as  $(v_i)_{i \geq 0}$  in descending order of  $T_v$ , and the edges  $e \in E(G)$  be ordered as  $(e_i)_{i \geq 0}$  in descending order of  $T_e$ . Then let  $H_V, H_E$  be the maximal prefixes of  $(v_i)_{i \geq 0}, (e_i)_{i \geq 0}$  such that  $\sum_{v \in H_V} T_v, \sum_{e \in H_E} T_e \leq \epsilon T$ . We define  $\Delta_{V, \epsilon}(G) = \max_{v \notin H_V} T_v, \Delta_{E, \epsilon}(G) = \max_{e \notin H_E} T_e$ .

When the graph meant is unambiguous, we will omit the parameter  $G$ .

**Lemma 35.** Let  $f$  be a triangle-dependent sampling algorithm that solves triangle counting for a graph  $G$  with  $(\epsilon, \delta)$  error. For any set  $A$  such that  $A \subseteq V(G) \cup E(G)$ , if the number of triangles in  $G$  involving edges or vertices from  $A$  is  $\geq 2\epsilon T$ , the probability that  $f$  samples no triangles involving edges or vertices from  $A$  must be  $\leq 2\delta$ .

*Proof.* Suppose this did not hold for some  $A$ . Then we can give the algorithm as input  $G$  with probability  $\frac{1}{2}$  and  $G \setminus A$  with probability  $\frac{1}{2}$ . Suppose  $f$  samples no triangles involving  $A$ . Then, as  $f$  is triangle-dependent, the output of the algorithm is independent of which of the two inputs the algorithm received. So there is at most a  $\frac{1}{2}$  chance the algorithm will return a value within  $\epsilon T(G)$  of its input. So this occurs with probability at most  $2\delta$ .  $\square$

We can now prove several lower bounds for such algorithms. Throughout, we will use the fact that for any sampling algorithm that is resilient to permutations of the input, the probability of sampling an edge can be taken to the same for all edges, or copies of any other fixed subgraph, without loss of generality. We will use  $m$  to refer to  $|E(G)|$  throughout.

**Lemma 36.** *Let  $f$  be a triangle-dependent sampling algorithm that solves triangle-counting for a graph  $G$  with probability  $\frac{9}{10}$  and error  $\epsilon$ . Then  $f$  samples  $\Omega\left(\frac{m}{T^{2/3}}\right)$  edges in expectation.*

*Proof.* This is a direct consequence of Lemma 32.  $\square$

**Lemma 37.** *Let  $f$  be a triangle-dependent sampling algorithm that solves triangle-counting a graph  $G$  with probability  $\frac{9}{10}$  and error  $\epsilon$ . Then  $f$  samples  $\Omega\left(m\frac{\sqrt{\Delta_{V,2\epsilon}}}{T}\right)$  edges in expectation.*

*Proof.* Our proof will make use of the fact that  $f$  must sample at least one triangle from the  $2\epsilon T$  triangles at the heaviest vertices. We will also use the fact that, as  $f$  needs to work for an arbitrary relabelling of the vertices, and because our sampling strategy is not allowed to vary based on which edges we've already seen, our sampling strategy will be the same for each vertex.

Let  $n$  be the number of vertices in  $G$ . Let  $X_v$  denote the distribution on how many of the  $n-1$  edges that could be incident to  $v$  will be sampled by  $f$  if they are in  $G$ . As we require our algorithm to work on all permutations of the vertices,  $X_v$  can be taken to be identically distributed for each  $v \in V(G)$ . Then, conditioned on  $X_v = i$ ,  $\Theta(i^2)$  of the  $\Theta(n^2)$  potential wedges next to  $v$  will be sampled if present. So the probability of sampling any given triangle that uses  $v$  is less than the chance of sampling its two edges incident to  $v$ , which is in turn  $\Theta(i^2/n^2)$ . Hence the chance of sampling a triangle at  $v$  is at most  $\Theta\left(\frac{i^2 T_v}{n^2}\right)$ .

Let the vertices  $v$  of  $G$  be ordered as  $(v_i)_{i \geq 0}$  in descending order of  $T_v$ . Let  $A$  be the minimal prefix of  $(v_i)_{i \geq 0}$  such that  $\sum_{v \in A} T_v \geq 2\epsilon T$ . Then as  $f$  is triangle-dependent, it must sample at least one triangle involving a vertex in  $A$  with probability at least  $\frac{4}{5}$ .

Conditioned on  $X_v = i$ , the chance of sampling a triangle at  $v$  is at most  $\Theta\left(\frac{i^2 T_v}{n^2}\right)$ , reaching  $\Theta(1)$  when  $\frac{i}{n} = \Theta(\sqrt{T_v}^{-1})$ . So, letting  $Y$  be the number of vertices  $v \in A$  such that at least one triangle is sampled at  $v$ ,  $\mathbb{E}[Y|X_v = i_v] \lesssim \sum_{v \in A} \min\left\{1, \frac{i_v^2 T_v}{n^2}\right\}$ , and so

$$\mathbb{E}[Y] \lesssim \sum_{v \in A} \mathbb{E}\left[\min\left\{1, \frac{X_v^2 T_v}{n^2}\right\}\right]$$

As the  $X_v$  are identically distributed,  $\mathbb{E}[X_v^2]$  is the same for all  $v$ . Given the constraints that  $\sum_{v \in A} T_v = O(T)$  and  $\forall v \in A, T_v \geq \Delta_{V,\epsilon}$ , the sum above will be maximised when there are  $O(T/\Delta_{V,\epsilon})$  vertices in  $A$ , each with  $T_v = \Delta_{V,\epsilon}$ . (as this is the configuration that minimizes the truncation of the tails.) So,

$$\mathbb{E}[Y] \lesssim \frac{T}{\Delta_{V,\epsilon}} \mathbb{E}\left[\min\left\{1, \frac{X_v^2 \Delta_{V,\epsilon}}{n^2}\right\}\right]$$

Now, for  $\mathbb{E}[Y]$  to be  $\Omega(1)$ , we need

$$\mathbb{E}\left[\min\left\{1, \frac{X_v^2 \Delta_{V,\epsilon}}{n^2}\right\}\right] \gtrsim \frac{\Delta_{V,\epsilon}}{T}$$

This, in turn, requires that  $\mathbb{E}[X_v]$  be at least (up to constants)  $\frac{\Delta_{V,\epsilon}}{T} \frac{n}{\sqrt{\Delta_{V,\epsilon}}} = \frac{n\sqrt{\Delta_{V,\epsilon}}}{T}$ . (as we minimize  $\mathbb{E}[X_v]$  by setting  $X_v$  to  $\frac{n}{\sqrt{\Delta_{V,\epsilon}}}$  with probability  $\frac{\Delta_{V,\epsilon}}{T}$ , and 0 otherwise.)

As, at each vertex, there are  $n$  possible edges, this implies sampling edges with  $\Omega\left(\frac{\sqrt{\Delta_{V,\epsilon}}}{T}\right)$  probability, and therefore sampling

$$\Omega\left(m\frac{\sqrt{\Delta_{V,2\epsilon}}}{T}\right)$$

edges in expectation. As we need  $\mathbb{E}[Y] = \Omega(1)$  to sample at least one triangle with  $4/5$  probability, this completes our proof.  $\square$

**Lemma 38.** *Let  $f$  be a triangle-dependent sampling algorithm that solves triangle-counting for a graph  $G$  with probability  $\frac{9}{10}$  and error  $\epsilon$ . Then  $f$  samples  $\Omega\left(m\frac{\Delta_{E,2\epsilon}}{T}\right)$  edges in expectation.*

*Proof.* Let the edges  $e$  of  $G$  be ordered as  $(e_i)_{i \geq 0}$  in descending order of  $T_e$ . Let  $A$  be the minimal prefix of  $(e_i)_{i \geq 0}$  such that  $\sum_{e \in A} T_e \geq 2\epsilon T$ . Then as  $f$  is triangle-dependent, it must sample at least one triangle involving an edge in  $A$  with probability  $\frac{4}{5}$ .

So as  $\Delta_{E,\epsilon} \leq T_e$  for all  $e \in A$ ,  $|A| \lesssim \frac{T}{\Delta_{E,\epsilon}}$ . So the algorithm must sample at a rate  $\gtrsim \frac{\Delta_{E,\epsilon}}{T}$ .  $\square$

**Theorem 4** (Triangle-dependent sampling bound). *For any constant  $\epsilon$  and for any graph  $G$ ,*

$$\text{INSTOPT}(G, \epsilon, 1/10) = \Omega\left(m\left(\frac{1}{T^{2/3}} + \frac{\sqrt{\Delta_{V,2\epsilon}}}{T} + \frac{\Delta_{E,2\epsilon}}{T}\right)\right)$$

*in the setting of triangle-dependent sampling algorithms.*

*Proof.* By combining the above three bounds.  $\square$

## F Refining the Upper Bound

**Definition 16.** *For any graph  $G$ ,  $\epsilon > 0$ , let the vertices  $v \in V(G)$  be ordered as  $(v_i)_{i \geq 0}$  in descending order of  $T_v$ , and the edges  $e \in E(G)$  be ordered as  $(e_i)_{i \geq 0}$  in descending order of  $T_e$ . Then let  $H_V, H_E$  be the maximal prefixes of  $(v_i)_{i \geq 0}, (e_i)_{i \geq 0}$  such that  $\sum_{v \in H_V} T_v, \sum_{e \in H_E} T_e \leq \epsilon T$ . We define  $\Delta_{V,\epsilon}(G) = \max_{v \notin H_V} T_v, \Delta_{E,\epsilon}(G) = \max_{e \notin H_E} T_e$ .*

*When the graph meant is unambiguous, we will omit the parameter  $G$ .*

**Theorem 17** (Refined triangle estimation upper bound). *If  $\tilde{T} \leq T$ , and  $\epsilon > 0$ , our algorithm obtains an  $(\epsilon, \delta)$  approximation to  $T$  while keeping*

$$O\left(\frac{m \log \frac{1}{\delta}}{\epsilon^2} \left(\frac{1}{\tilde{T}^{2/3}} + \frac{\sqrt{\Delta_{V,\epsilon/24}} \log \tilde{T}}{\tilde{T}} + \frac{\Delta_{E,\epsilon/24} \log \tilde{T}}{\tilde{T}}\right)\right)$$

*edges. If  $\tilde{T} > T$ , the algorithm either determines that  $\tilde{T} > T$  or obtains a  $(1 \pm \epsilon)$  approximation to  $T$  with probability  $1 - \delta$ .*

*Proof.* As in the proof of Theorem 1, the variance of the output  $\bar{T}$  of a single iteration of the algorithm will be  $O(T\tilde{T})$  if  $k = \min\left\{\tilde{T}^{2/3}, \frac{\tilde{T}^{3/2}}{\sqrt{\Delta_V}}, \frac{\tilde{T}^2}{\Delta_E}\right\}$ .

Now define the graph  $G'$  as follows: order the vertices  $v \in V(G)$  as  $(v_i)_{i \geq 0}$  in descending order of  $T_v$ , and let  $V'$  be the minimal postfix of  $(v_i)_{i \geq 0}$  such that  $\sum_{v \in V(G) \setminus V'} T_v \leq \frac{\epsilon}{24} T$ . Order the edges  $e \in E(G)$  as  $(e_i)_{i \geq 0}$  in descending order of  $T_e$ , and let  $E'$  be the minimal postfix of  $(e_i)_{i \geq 0}$  such that  $\sum_{e \in E(G) \setminus E'} T_e \leq \frac{\epsilon}{24} T$ .  $G'$  is  $(V', E' \cap (V' \times V'))$ .

Let  $\bar{T}_\epsilon$  be  $\bar{T}$ , less any contributions *in the second pass* from triangles not contained in  $G'$ . We can then compare this to  $\bar{T}(G')$ , the result the algorithm would have given if run with  $G'$  as input. The only difference between  $\bar{T}_\epsilon$  and  $\bar{T}(G')$  will come from triangles counted in the *first* pass. As the

expectation of  $\bar{T}(G')$  is independent of the output of the first pass,  $\mathbb{E}[\bar{T}_\epsilon] = \mathbb{E}[\bar{T}(G')] = T(G') = (1 - \epsilon')T$  for some  $\epsilon' \in [0, \epsilon/12]$ . Furthermore, as counting additional triangles in the second pass either has no effect or increases the estimate of  $T$ ,  $\bar{T} - \bar{T}_\epsilon$  will always be  $\geq 0$ .

We now seek to bound the variance of  $\bar{T}_\epsilon$ . We do so by comparison to the variance of  $\bar{T}(G')$ . As the execution of the two algorithms differs only in the first pass, and therefore in the values of the “triangle degree estimates” for each vertex, we need consider only how this impacts lemmas 20, 21, and 22, the three results about the distributions of these estimates on which the variance depends. This will allow us to show that the proof of our bound on  $\text{Var}(\bar{T}(G'))$  is, within a constant factor, a bound on  $\text{Var}(\bar{T}_\epsilon)$ .

For clarity, we will use  $\mathcal{T}_v$  to refer to the estimate for the vertex  $v$  in the calculation of  $\bar{T}_\epsilon$  (as it is identical to the estimate calculated for  $\bar{T}$ ), and  $\mathcal{T}'_v$  for the estimate used in the calculation of  $\bar{T}$ . As  $\mathcal{T}_v$  is monotonic decreasing in the number of first-pass triangles counted at  $v$ ,  $\mathcal{T}_v \leq \mathcal{T}'_v$  for all  $v$ .

We note that  $\mathcal{T}_v \leq \mathcal{T}'_v$  for all  $v$ , as  $\mathcal{T}_v$  is monotonic decreasing in the number of first-pass triangles counted at  $v$ . Furthermore, if  $\mathcal{T}_v$  is LIGHT, so is  $\mathcal{T}'_v$ . Therefore, Lemmas 20 and 22 will hold without further work. Lemma 21, however, gives us a bound on  $\mathcal{T}'_v$  not being LIGHT, in terms of the parameters of  $G'$ , which does not directly transfer to a bound on  $\mathcal{T}_v$ .

**Lemma 21.**  $\forall v \in V, \mathbb{P}[\mathcal{T}_v \neq \text{LIGHT}] \lesssim \frac{\sqrt{k}\mathcal{T}_v}{\bar{T}}$ .

The lemma does, however, give us a bound on  $\mathcal{T}_v$  in terms of the parameters  $T_v(G)$  and  $\tilde{T}$ , rather than  $T_v(G')$ . The only place this lemma is used in bounding the  $l = 1, u \neq v$  case in Lemma 26, and so substituting in this bound instead of the one we would normally benefit from, we replace the  $T_E(G')k\frac{T(G')}{\bar{T}}$  term in the variance with

$$\begin{aligned} \sqrt{k}T_E(G') \sum_{v \in V(G')} \mathbb{P}[\mathcal{T}_v \neq \text{LIGHT}] &\lesssim \sqrt{k}T_E(G') \sum_{v \in V(G')} \frac{\sqrt{k}\mathcal{T}_v(G)}{\tilde{T}} \\ &= T_E(G')k\frac{T(G)}{\tilde{T}} \\ &\lesssim T_E(G')k\frac{T(G')}{\tilde{T}} \end{aligned}$$

As  $T(G') > (1 - \epsilon/12T(G))$ . Therefore, we lose at most a constant factor in going from our bound on the variance of  $\bar{T}(G')$  to bounding the variance of  $\bar{T}_\epsilon$ . So  $\text{Var}(\bar{T}_\epsilon) \lesssim T\tilde{T}$  if  $k = \min\left\{\tilde{T}^{2/3}, \frac{\tilde{T}^{\frac{3}{2}}}{\sqrt{\Delta_{V,\epsilon/24}}}, \frac{\tilde{T}^2}{\Delta_{E,\epsilon/24}}\right\}$ . (as  $\Delta_V(G') \leq \Delta_{V,\epsilon}(G)$  and  $\Delta_E(G') \leq \Delta_{E,\epsilon}(G)$ )

We can therefore split  $\bar{T}$  into two (not independent) random variables,  $\bar{T}_\epsilon$  and  $T - \bar{T}_\epsilon$ , with  $\mathbb{E}[\bar{T}_\epsilon] = (1 - \epsilon')T$ ,  $\mathbb{E}[T - \bar{T}_\epsilon] = \epsilon'T$ , and  $\text{Var}(\bar{T}_\epsilon) \lesssim T\tilde{T}$ . Now, we can repeat the algorithm  $c$  times, averaging the results as  $\bar{T}^{(c)}$ . We then have  $\bar{T}^{(c)} = \bar{T}_\epsilon^{(c)} + (T - \bar{T}_\epsilon^{(c)})$ , where  $\bar{T}_\epsilon^{(c)}$  is the average of the  $c$  instances of  $\bar{T}_\epsilon$ . So  $\mathbb{E}[\bar{T}_\epsilon^{(c)}] = (1 - \epsilon')T$ ,  $\mathbb{E}[\bar{T}^{(c)} - \bar{T}_\epsilon^{(c)}] = \epsilon'T$ , and  $\text{Var}(\bar{T}_\epsilon^{(c)}) \lesssim \frac{1}{c}T\tilde{T}$ . So we can choose  $c = O(\frac{1}{\epsilon})$  to get  $\text{Var}(\bar{T}_\epsilon^{(c)}) \leq \frac{\epsilon}{4}T\tilde{T}$ .

Then, by Chebyshev’s inequality,  $\bar{T}_\epsilon^{(c)}$  is within  $\frac{\epsilon}{2}\sqrt{T\tilde{T}}$  of  $(1 - \epsilon')T$  with probability  $\frac{3}{4}$ , and by Markov’s inequality (and the fact that it is non-negative),  $\bar{T}^{(c)} - \bar{T}_\epsilon^{(c)}$  is  $< 6\epsilon'T \leq \frac{\epsilon}{2}T$  with probability  $\frac{5}{6}$ . So by taking a union bound, both of those hold with probability  $\frac{2}{3}$ . So if  $\tilde{T} \leq T$ ,  $\bar{T}$  will be  $(1 \pm \epsilon)T$  with probability  $\frac{2}{3}$ , and if  $\tilde{T} \geq T$ ,  $\bar{T}$  will be within  $\epsilon\tilde{T}$  of  $T$  with probability  $\frac{2}{3}$ . (so, by responding that  $T < \tilde{T}$  if  $\bar{T} < (1 - \epsilon)\tilde{T}$ , and returning  $\bar{T}$  otherwise, it will either determine  $T < \tilde{T}$  correctly or return a  $(1 \pm \epsilon)$  approximation to  $T$ )

Then, by repeating this process  $O(\log \frac{1}{\delta})$  times and taking the median, we can guarantee that the above will hold with probability  $1 - \delta$  instead, giving us our final result.  $\square$

## G Counting Constant Size Subgraphs

We now consider the problem of counting subgraphs of a constant size in our graph stream. Let  $A$  be a fixed subgraph, with  $s = |A|$ . We will attempt to estimate  $M$ , the number of subgraphs of  $G$  that are isomorphic to  $A$ .

We will use the same algorithm as for triangles, slightly generalized.  $\widetilde{M}$  will be a lower bound on  $M$ .  $\forall S \subseteq V(G)$ ,  $\mu(S)$  is the set of copies of  $A$  in  $G$  that contain  $S$ ,  $M_S = |\mu(S)|$ , and  $\forall i \in [s]$ ,  $C_i = \sum_{S \subseteq V(G), |S|=i} T_S^2$ . With  $C_i^+$  as an upper bound on  $C_i$ ,  $\omega$  is then defined, analogously to the triangle case, as  $\min \left\{ \frac{(\widetilde{M})^2}{C_1^+}, \sqrt{k} \right\}$ .

### G.1 First Pass

For  $i = 1, 2$ : Let  $r_{D,i} : E \rightarrow [0, 1]$  be a uniformly random hash function.

Let  $\mathcal{M}_i : V \rightarrow \{0, 1\}$  be a random hash function such that  $\forall v \in V$ :

$$\mathcal{M}_i(v) = \begin{cases} 1 & \text{With probability } \frac{1}{\sqrt{k}}. \\ 0 & \text{Otherwise.} \end{cases}$$

For each  $v$  in  $V(G)$ ,  $a \in \mu(\{v\})$ , and  $h \in \left\{ 0, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil \right\}$ , define  $X_{v,a,i}^{(h)}$  to be 1 if the following hold:

$$\begin{aligned} \forall u \in V(a) \setminus \{v\}, \mathcal{M}_i(u) = 1 \\ \forall \{u|vu \in E(a)\}, r_{D,i}(vu) < \frac{\omega 2^h}{\sqrt{k}} \end{aligned}$$

And 0 otherwise. We define  $x_{v,a,i}^{(h)} = \mathbb{E} \left[ X_{v,a,i}^{(h)} \right]$ , so if  $t = |\{vu|uv \in E(a)\}|$ ,  $x_{v,a,i}^{(h)} = \frac{\omega^t 2^{ht}}{k^{\frac{s+t-1}{2}}}$ .

Then let  $X_{v,i}^{(h)} = \sum_{a \cong A, v \in V(a)} \frac{X_{v,a,i}^{(h)}}{x_{v,a,i}^{(h)}}$ .

We then define  $H_{v,i} = \left\{ h \in \left\{ 0, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil \right\} \mid X_{v,i}^{(h)} \geq \frac{\widetilde{M}}{\omega 2^h} \right\}$ .  $\mathcal{M}_{v,i}$  is then defined as follows:

$$\mathcal{M}_{v,i} = \begin{cases} \min H_{v,i} & \text{If } H_{v,i} \neq \emptyset. \\ \text{LIGHT} & \text{Otherwise.} \end{cases}$$

Then, for each  $v \in V(G)$ , we define  $\mathcal{M}_v$  to be LIGHT if  $\mathcal{M}_{v,1}$  and  $\mathcal{M}_{v,2}$  are LIGHT, and otherwise to be the smallest numerical value amongst  $\mathcal{M}_{v,1}, \mathcal{M}_{v,2}$ .

### G.2 Second Pass

#### G.2.1 Splitting the Graph

Let  $V_L = \{v \in V \mid \mathcal{M}_v = \text{LIGHT}\}$ , and let  $G_L$  be the subgraph induced by  $V_L$ . Then we define  $M_L$  as the number of copies of  $A$  in  $G_L$ , and  $M_H = M - M_L$ .

We will compute estimates  $\overline{M}_L, \overline{M}_H$  of  $M_L, M_H$ , and estimate  $M$  as  $\overline{M} = \overline{M}_L + \overline{M}_H$ .

### G.3 Estimating $M_L$

We will estimate  $M_L$  by sampling the vertices of  $V_L$  with probability  $\frac{1}{\sqrt{k}}$  each and calculating the number of copies of  $A$  in the resulting graph.

Let  $c : V \rightarrow \{0, 1\}$  be a random hash function such that  $\forall v \in V$ :

$$c(v) = \begin{cases} 1 & \text{With probability } \frac{1}{\sqrt{k}}. \\ 0 & \text{Otherwise.} \end{cases}$$

Let  $V'_L = \{v \in V_L | c(v) = 1\}$ , and let  $G'_L$  be the subgraph of  $G$  induced by  $V'_L$ . Then  $\overline{M}_L$  is the number of copies of  $A$  in  $G'_L$ , multiplied by  $k^{\frac{s}{2}}$ .

### G.4 Estimating $M_H$

We will estimate  $M_H$  by considering every vertex in  $V \setminus V_L$  separately. We will achieve this by sampling vertices  $v$  with probability  $2^{-\mathcal{M}_v}$ , and then sampling edges incident to  $v$  with probability proportional to  $2^{\mathcal{M}_v}$ .

Let  $h : V \rightarrow \{0, 1\}$  be a random hash function such that  $\forall v \in V$ :

$$h(v) = \begin{cases} i & \text{With probability } \frac{1}{\omega 2^i} \text{ for each } i \in \{0, \dots, \lceil \log \frac{\sqrt{k}}{\omega} \rceil\}. \\ -\infty & \text{Otherwise.} \end{cases}$$

And let  $r_C : V \rightarrow [0, 1]$  be a uniformly random hash function.

Then, for each  $v \in V_H = V \setminus V_L$ , we allow  $v$  to contribute to  $\overline{M}_L$  iff  $h(v) = \mathcal{M}_v$ . If it does, we calculate its contribution  $\overline{M}_v$  as follows:

We count a subgraph  $a \subseteq$  such that  $a \cong A$  and  $v \in V(a)$  iff:

$$\begin{aligned} & \forall u \in V(a) \setminus \{v\}, c(u) = 1 \\ & \forall e \in \{u|vu \in E(a)\}, r_C(e) < \frac{\omega 2^h}{\sqrt{k}} \end{aligned}$$

Then, for each such  $a$ , let  $t = |\{vu | u \in E(a)\}|$ . We add  $\frac{k^{\frac{s+t-1}{2}}}{\omega^{t-1} 2^{\mathcal{M}_v(t-1)}} \times \frac{1}{|\{x \in V(a) | \mathcal{M}_x \neq \text{LIGHT}\}|}$  to  $\overline{M}_v$ . (with the second term then compensating for the fact that a motif could potentially be counted at multiple different vertices)

We then define  $\overline{M}_H = \sum_{v \in V_H, h(v) = \mathcal{M}_v} \overline{M}_v$ .

### G.5 Final Output

We then output our estimate  $\overline{M}$  of  $M$  as  $\overline{M} = \overline{M}_L + \overline{M}_H$ .

## H Analysis of Generalized Algorithm

### H.1 First Pass

#### H.1.1 Relation of $\mathcal{M}_v$ to $M_v$

**Lemma 39.** *For any  $v \in [n]$ , let  $X_v^{(h)}$  be as defined previously. Then,  $\mathbb{E} [X_v^{(h)}] = M_v$  and  $\text{Var} (X_v^{(h)}) \leq \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left(\frac{k}{\omega 2^h}\right)^{l-1}$ .*

*Proof.*  $X_v^{(h)} = \sum_{a \in \mu(\{v\})} \frac{X_{v,a}^{(h)}}{\mathbb{E}[X_{v,a}^{(h)}]}$ , so  $\mathbb{E}[X_v^{(h)}] = M_v$ .

Then, to bound the variance, we start by bounding  $\mathbb{E}\left[\left(X_v^{(h)}\right)^2\right] = \sum_{a_1, a_2 \in \mu(\{v\})} \frac{\mathbb{E}[X_{v,a_1}^{(h)} X_{v,a_2}^{(h)}]}{x_{v,a_1}^{(h)} x_{v,a_2}^{(h)}}$ . We will consider the contribution of these terms for each possible value of  $l = |V(a_1) \cap V(a_2)|$ . (noting that, for any such subgraph  $a$ , if  $t = |\{vu|u \in E(a)\}|$ ,  $t \leq s-1$ ):

Letting  $t_i = |\{vu|vu \in E(a_i)\}|$  for  $i = 1, 2$ ,  $x_{v,a_i}^{(h)} = \frac{\omega^{t_i} 2^{ht_i}}{k^{\frac{s+t_i-1}{2}}}$ .

Then, for  $X_{v,a_1}^{(h)} X_{v,a_2}^{(h)} = 1$  to hold, we need the  $2s-l-1$  vertices  $u \in V(a_1) \cup V(a_2) \setminus \{v\}$  to have  $d(u) = 1$ , and the  $\geq t_1 + t_2 - (l-1)$  edges  $e \in \{vu|vu \in E(a_1)\} \cup \{vu|vu \in E(a_2)\}$  to have  $r_D(e) < \frac{\omega 2^h}{\sqrt{k}}$ . So  $\mathbb{E}[X_{v,a_1}^{(h)} X_{v,a_2}^{(h)}] \leq \frac{(\omega 2^h)^{t_1+t_2+1-l}}{k^{(s-l)+\frac{t_1+t_2}{2}}}$ .

So  $\frac{\mathbb{E}[X_{v,a_1}^{(h)} X_{v,a_2}^{(h)}]}{x_{v,a_1}^{(h)} x_{v,a_2}^{(h)}} \leq \left(\frac{k}{\omega 2^h}\right)^{l-1}$ . There are  $\leq \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2$  such pairs for each  $l \in \{2, s\}$ , and so the total contribution to the sum is  $\leq \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left(\frac{k}{\omega 2^h}\right)^{l-1}$ .

This gives us:

$$\begin{aligned} \text{Var}\left(X_v^{(h)}\right) &= \mathbb{E}\left[\left(X_v^{(h)}\right)^2\right] - \mathbb{E}\left[X_v^{(h)}\right]^2 \\ &\leq \sum_{i=1}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left(\frac{k}{\omega 2^h}\right)^{l-1} - M_v^2 \\ &= \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left(\frac{k}{\omega 2^h}\right)^{l-1} \end{aligned}$$

□

**Lemma 40.** *If  $M_v \geq 2\frac{\widetilde{M}}{\sqrt{k}}$ , then  $\mathbb{P}[\mathcal{M}_v = \text{LIGHT}] \lesssim \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 k^{\frac{l-1}{2}}$ .*

*Proof.* For  $\mathcal{M}_v = \text{LIGHT}$  to hold, we need that  $X_v^{(h)} < \frac{\widetilde{M}}{\omega 2^h}$  for all  $h \in \left[\left[\log \frac{\sqrt{k}}{\omega}\right]\right]$ , and so in particular  $X_v^{\left(\left[\log \frac{\sqrt{k}}{\omega}\right]\right)} < \frac{\widetilde{M}}{\sqrt{k}}$ .

By Lemma 39,  $\mathbb{E}\left[X_v^{\left(\left[\log \frac{\sqrt{k}}{\omega}\right]\right)}\right] = M_v \geq 2\frac{\widetilde{M}}{\sqrt{k}}$ , and

$$\text{Var}\left(X_v^{\left(\left[\log \frac{\sqrt{k}}{\omega}\right]\right)}\right) \lesssim \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 k^{\frac{l-1}{2}}$$

. So by Chebyshev's inequality, this holds with probability  $\lesssim \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 k^{\frac{l-1}{2}}$ . □

**Lemma 41.**  $\forall v \in V, \mathbb{P}[\mathcal{M}_v \neq \text{LIGHT}] \lesssim \frac{\sqrt{k} M_v}{M}$ .

*Proof.* For  $\mathcal{M}_v \neq \text{LIGHT}$  to hold, we need there to be at least one  $h \in \left[\left[\log \frac{\sqrt{k}}{\omega}\right]\right]$  such that  $X_v^{(h)} \geq \frac{\widetilde{M}}{\omega 2^h}$ . By Lemma 39,  $\mathbb{E}[X_v^{(h)}] = M_v$ , so by Markov's inequality this occurs with probability

$\leq \frac{M_v}{M} \omega 2^h$ . So the probability that it holds for any  $h$  is  $\leq \sum_{h=\lceil \log \frac{\sqrt{k}}{\omega} \rceil}^1 \frac{M_v}{M} \omega 2^h \lesssim \sum_{i=0}^{\infty} \frac{M_v \sqrt{k}}{M} 2^{-i} \lesssim \frac{\sqrt{k} M_v}{M}$ .  $\square$

**Lemma 42.**  $\forall v \in V, l \geq 1, \lceil \log \frac{\widetilde{M}}{\omega M_v} \rceil \geq 1,$

$$\mathbb{P} \left[ \mathcal{M}_v \geq \max \left\{ \left\lceil \log \frac{\widetilde{M}}{\omega M_v} \right\rceil + j, 0 \right\} \right] \lesssim \left( \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{k M_v}{2^j \widetilde{M}} \right)^{l-1} \right)^2$$

*Proof.* For  $\mathcal{M}_{v,i} \geq \max \left\{ \left\lceil \log \frac{\widetilde{M}}{\omega M_v} \right\rceil + j, 0 \right\}$  to hold,  $X_v^{\max \left\{ \left\lceil \log \frac{\widetilde{M}}{\omega M_v} \right\rceil + j - 1, 1 \right\}} < \frac{M_v}{2}$  must hold. By Lemma 39 and Chebyshev's inequality, this happens with probability

$$\lesssim \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{k M_v}{2^j \widetilde{M}} \right)^{l-1}$$

. So it holds for both  $i = 1, 2$ , and therefore for  $\mathcal{M}_v$ , with probability

$$\lesssim \left( \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{k M_v}{2^j \widetilde{M}} \right)^{l-1} \right)^2$$

$\square$

## H.2 Second Pass

### H.2.1 $\overline{M}_L$

**Lemma 43.**  $\mathbb{E} [\overline{M}_L | M_L] = M_L.$

*Proof.*  $M_L$  is the number of copies of  $A$  in the random graph  $G_L$ , and  $\overline{M}_L$  is  $k^{\frac{s}{2}}$  times the number of triangles in the random graph  $G'_L$ . Any subgraph  $a \subseteq G$  is in  $G'_L$  iff  $\forall v \in V(a), c(v) = 1$ , which occurs with probability  $\frac{1}{k^{\frac{s}{2}}}$  for subgraphs  $a$  such that  $a \cong A$ . So the expected number of copies in  $G'_L$  is  $\frac{M_L}{k^{\frac{s}{2}}}$ , and so  $\mathbb{E} [\overline{M}_L | M_L] = M_L$ .  $\square$

**Lemma 44.**  $\text{Var}(\overline{M}_L) \lesssim M^2 + \sum_{l=2}^s C_l k^{\frac{l}{2}}.$

*Proof.* Let  $M'_L$  be the number of copies of  $A$  in  $G'_L$ . (so  $\overline{M}_L = k^{\frac{s}{2}} M'_L$ ) A subgraph  $a \cong A$  in  $G$  is in  $G'_L$  iff  $\forall v \in V(a), c(v) = 1$  and  $\mathcal{M}_v = \text{LIGHT}$ . We proceed by bounding  $\mathbb{E} [M'^2_L]$ , by considering, for any pair of subgraphs  $a_1, a_2 \cong A$ , the probability they are both in  $G'_L$ . We will bound this for each value of  $l = |V(a_1) \cap V(a_2)|$ , treating  $l = 1$  as a special case.

$l \neq 1$ :  $\mathbb{P} [\forall v \in V(a), c(v) = 1] = k^{-\frac{|V(a_1) \cap V(a_2)|}{2}} = k^{-\frac{2s-l}{2}}$ . So the total contribution to  $\mathbb{E} [M'^2_L]$  from these pairs is  $\leq C_l k^{\frac{l-2s}{2}}$ .

$l = 1$ : Let  $v$  be the unique vertex in  $V(a_1) \cap V(a_2)$ . Then in this case, we will also make use of the fact that  $\mathcal{M}_v = \text{LIGHT}$  must hold. As above,  $\mathbb{P} [\forall v \in V(a), c(v) = 1] = k^{-\frac{2s-1}{2}}$ .



So for vertices  $v$  s.t.  $M_v \leq \frac{2\widetilde{M}}{\sqrt{k}}$ , we can bound the contribution to  $\mathbb{E}[M_L'^2]$  from these cases by  $k^{-\frac{2s-1}{2}} \sum_{v \in V(G)} M_v^2 \lesssim k^{-\frac{2s-1}{2}} \frac{(\widetilde{M})^2}{\sqrt{k}} = k^{-s} (\widetilde{M})^2 \leq k^{-s} M^2$ .

We can then consider vertices  $v$  s.t.  $M_v > \frac{2\widetilde{M}}{\sqrt{k}}$ . By Lemma 40,

$$\mathbb{P}[\mathcal{M}_v = \text{LIGHT}] \lesssim \frac{1}{M_v^2} \sum_{i=2}^s \sum_{S \subseteq V(G), |S|=i, v \in S} M_S^2 k^{\frac{i-1}{2}}$$

. As  $\mathcal{M}_v$  is independent of the hash function,  $c$ , the probability of  $a_1$  and  $a_2$  both being in  $G'_L$  is  $\lesssim \frac{1}{M_v^2} \sum_{i=2}^s \sum_{S \subseteq V(G), |S|=i, v \in S} M_S^2 k^{\frac{i-2s}{2}}$ , and so the total contribution to  $\mathbb{E}[M_L'^2]$  is  $\lesssim \sum_{v \in V(G)} \sum_{i=2}^s \sum_{S \subseteq V(G), |S|=i, v \in S} M_S^2 k^{\frac{i-2s}{2}} \lesssim \sum_{i=2}^s C_i k^{\frac{i-2s}{2}}$ .

So by summing these cases together:

$$\mathbb{E}[M_L'^2] \lesssim k^s M^2 + \sum_{l=2}^s C_l k^{\frac{l-2s}{2}}$$

Which then gives us our variance bound:

$$\begin{aligned} \text{Var}(\overline{M}_L) &\leq \mathbb{E}[M_L^2] \\ &\leq k^s \mathbb{E}[M_L'^2] \\ &\lesssim M^2 + \sum_{l=2}^s C_l k^{\frac{l}{2}} \end{aligned}$$

□

## H.2.2 $\overline{M}_H$

**Lemma 45.** *If  $\widetilde{M} \leq M$ ,  $\mathbb{E}[\overline{M}_H | M_H] = M_H$ .*

*Proof.*  $\overline{M}_H = \sum_{v \in V_H} \overline{M}_v$ , and  $M_H$  is the number of copies of  $A$  in  $G$  that are not in  $G_L$ . If a subgraph in  $G$  is in  $G_L$ , it can never contribute to a  $T_v$ , as every vertex  $u$  it uses will have  $\mathcal{M}_u = \text{LIGHT}$ . If a subgraph  $a \cong A$  in  $G$  is not in  $G_L$ , it has  $l \in [s]$  vertices  $u$  s.t.  $\mathcal{M}_u \neq \text{LIGHT}$ . It will therefore have  $l$  vertices  $v$  such that it can contribute to  $M_v$ .

At each of those vertices  $v$ , a subgraph  $a \cong A$  with  $t = |\{u | vu \in E(a)\}|$  will contribute  $\frac{k^{\frac{s+t-1}{2}}}{\omega^{t-1} 2^{\mathcal{M}_v(t-1)}} \frac{1}{l}$  to  $\overline{M}_v$  iff:

$$\begin{aligned} h(v) &= \mathcal{M}_v \\ \forall u \in V(a) \setminus \{v\}, c(u) &= 1 \\ \forall e \in \{u | vu \in E(a)\}, r_C(e) &< \frac{\omega 2^{\mathcal{M}_v}}{\sqrt{k}} \end{aligned}$$

This happens with probability  $\frac{(\omega 2^{\mathcal{M}_v})^t}{\omega 2^{\mathcal{M}_v} k^{\frac{s+t-1}{2}}}$ , so the expected contribution of  $a$  to  $\overline{M}_v$  is  $\frac{1}{l}$ . Therefore, as there are  $l$  vertices where  $a$  can contribute, its expected contribution to  $\overline{M}_H$  is 1.

Therefore,  $\mathbb{E}[\overline{M}_H | M_H]$  is precisely the number of copies of  $A$  in  $G$  that are not in  $G_L$ , which is  $M_H$ . □

**Lemma 46.** *If  $\widetilde{M} \leq M$ ,  $\text{Var}(\overline{M}_H) \lesssim M^2 + \sum_{l=2}^s C_l \frac{k^{l-1}}{\omega^{l-2}}$*

*Proof.* We now care, for any subgraph  $a \cong A$ , which vertex we are counting it at (and so which  $M_v$  it may contribute to. We will therefore use  $a^v$  to denote the subgraph  $a$ , counted at the vertex  $v$ . With  $t = |\{u|vu \in E(a)\}|$ , we then define  $Y_{a^v}$  as  $\frac{k^{\frac{s+t-1}{2}}}{\omega^{t-1}2^{\mathcal{M}_v(t-1)}}$  if  $a$  is counted at  $v$  and 0 otherwise.

Then,  $\overline{M}_v \leq \sum_{a \subseteq G, a \cong A, v \in V(a)} Y_{a^v}$ . So with  $Y = \sum_{a, v, a \cong A, v \in V(a)} Y_{a^v}$ ,  $\mathbb{E}[\overline{M}^2] \leq \mathbb{E}[Y^2]$ . We will bound  $\mathbb{E}[Y^2]$  by bounding  $\mathbb{E}[Y_{a_1^u} Y_{a_2^v}]$  for each pair  $a_1^u, a_2^v$ .

We will bound these with respect to  $l = |V(a_1) \cap V(a_2)|$ , treating  $l = 1$  as a special case, and treating  $u = v$  and  $u \neq v$  separately. In each case, let  $t_1 = |\{uw|uw \in E(a_1)\}|$ ,  $t_2 = |\{vw|vw \in E(a_2)\}|$ .

$l \neq 1, u \neq v$ : We need  $\mathcal{M}_u = h(u)$ ,  $\mathcal{M}_v = h(v)$ , which happens with probability  $\omega^{-2}2^{-\mathcal{M}_u - \mathcal{M}_v}$ . We need  $\forall w \in V(a_1) \setminus \{u\} \cup V(a_2) \setminus \{v\}, c(w) = 1$ . There are at least  $2s - l - 2$  vertices in this set, so this happens with probability  $\leq k^{-s+1+\frac{l}{2}}$ . We need  $\forall e \in \{uw|uw \in E(a_1)\} r_C(e) < \frac{\omega 2^{\mathcal{M}_u}}{\sqrt{k}}$  and  $\forall e \in \{vw|vw \in E(a_2)\}, r_C(e) < \frac{\omega 2^{\mathcal{M}_v}}{\sqrt{k}}$ . These sets can overlap in at most one edge  $(uv)$ , so this happens with probability  $\leq \left(\frac{\omega}{\sqrt{k}}\right)^{t_1+t_2-1} 2^{t_1 \mathcal{M}_u + t_2 \mathcal{M}_v - \max\{\mathcal{M}_u, \mathcal{M}_v\}}$ . Furthermore, the overlap in  $uv$  can only occur if  $u \in V(a_2)$  and  $v \in V(a_1)$ , and in this case we will also need  $c(u) = 1$  and  $c(v) = 1$ . So this either reduces the probability by a factor of  $\frac{\omega 2^{\max\{\mathcal{M}_u, \mathcal{M}_v\}}}{\sqrt{k}}$  or  $\frac{1}{k}$ , and so in either case by at least a factor of  $\frac{\omega 2^{\max\{\mathcal{M}_u, \mathcal{M}_v\}}}{\sqrt{k}}$ .

So as these three conditions are independent, and multiplying by the values  $Y_{a_1^u}, Y_{a_2^v}$  take when  $a_1^u, a_2^v$  are counted,  $\mathbb{E}[Y_{a_1^u} Y_{a_2^v}] \leq k^{\frac{l}{2}}$ . So the contribution to  $\mathbb{E}[Y^2]$  from such pairs is  $\lesssim C_l k^{\frac{l}{2}}$ .

$l \neq 1, u = v$ : We need  $\mathcal{M}_v = h(v)$ , which happens with probability  $\omega^{-1}2^{-\mathcal{M}_v}$ . We need  $\forall w \in V(a_1) \setminus \{u\} \cup V(a_2) \setminus \{v\}, c(w) = 1$ . There are at least  $2s - l - 1$  vertices in this set, so this happens with probability  $\leq k^{-s+\frac{l+1}{2}}$ . We need

$$\forall e \in \{uw|uw \in E(a_1)\} \cup \{vw|vw \in E(a_2)\}, r_C(e) < \frac{\omega 2^{\mathcal{M}_v}}{\sqrt{k}}$$

. As this set has at least  $\max\{t_1, t_2\}$  elements, this happens with probability  $\leq \left(\frac{\omega 2^{\mathcal{M}_v}}{\sqrt{k}}\right)^{\max\{t_1, t_2\}}$ .

So as these three conditions are independent, and multiplying by the values  $Y_{a_1^u}, Y_{a_2^v}$  take when  $a_1^u, a_2^v$  are counted,  $\mathbb{E}[Y_{a_1^u} Y_{a_2^v}] \leq k^{\frac{l+\min\{t_1, t_2\}-3}{2}} \omega^{1-\min\{t_1, t_2\}} 2^{1-\mathcal{M}_v(\min\{t_1, t_2\})}$ . So as  $\frac{\omega 2^{\mathcal{M}_v}}{\sqrt{k}} \leq 1$ , and  $t_1, t_2 \leq l - 1$ , this gives us  $\mathbb{E}[Y_{a_1^u} Y_{a_2^v}] \leq \frac{k^{l-1}}{\omega^{l-2} 2^{\mathcal{M}_v(l-2)}} \leq \frac{k^{l-1}}{\omega^{l-2}}$ . So the contribution to  $\mathbb{E}[Y^2]$  from such pairs is  $\lesssim C_l \frac{k^{l-1}}{\omega^{l-2}}$ .

$l = 1, u \neq v$ : As in the  $l \neq 1$  case,  $\mathbb{E}[Y_{a_1^u} Y_{a_2^v}] \leq k^{\frac{l}{2}} = \sqrt{k}$ . So we seek to bound the number of such pairs.

For any  $w \in V(G)$ , the number of such pairs intersecting at  $w$  is  $\leq \left(\sum_{v \in V(G), \mathcal{M}_v \neq \text{LIGHT}} M_{wv}\right)^2$ . Now let  $L = |\{v \in V(G) | \mathcal{M}_v \neq \text{LIGHT}\}|$ . Suppose  $L = r$ . By Cauchy-Schwartz, this means the number of such pairs intersecting at  $w$  is  $\leq r \sum_{v \in V(G)} M_{wv}^2$ . So by summing across all  $w$ , the total number of such pairs is  $\leq r C_2$ . So the contribution to the expectation conditioned

on  $L = r$  is  $\leq \sqrt{kr}C_2$ . As our bound on  $\mathbb{E}[Y_{a_1^u}Y_{a_2^v}]$  is independent of the values of the  $\mathcal{M}_v$ , we can then bound the unconditional contribution to  $\mathbb{E}[Y^2]$  by:

$$\begin{aligned}
\sum_r \sqrt{kr}C_2\mathbb{P}[L=r] &= \mathbb{E}[L]\sqrt{k}C_2 \\
&= \sqrt{k}C_2 \sum_{v \in V(G)} \mathbb{P}[\mathcal{M}_v = \text{LIGHT}] \\
&= \sqrt{k}C_2 \sum_{v \in V(G)} \frac{\sqrt{k}M_v}{\widetilde{M}} && \text{By Lemma 41.} \\
&= C_2k \frac{M}{\widetilde{M}}
\end{aligned}$$

$l = 1, u = v$ : As in the  $l \neq 1$  case,  $\mathbb{E}[Y_{a_1^u}Y_{a_2^v}] \leq \frac{k^{l-1}}{\omega^{l-2}2^{\mathcal{M}_v(l-2)}} = \omega 2^{\mathcal{M}_v}$ . So at any vertex  $v$ , the contribution to the expectation, conditioned on  $\mathcal{M}_v$ , is  $\leq M_v^2 \omega 2^{\mathcal{M}_v}$ .

We will consider three (exhaustive, but not necessarily mutually exclusive) cases:  $\mathcal{M}_v \leq \left\lceil \log \frac{\widetilde{M}}{\omega M_v} \right\rceil$ ,  $\mathcal{M}_v = 0$ , and  $\mathcal{M}_v = \max \left\{ \left\lceil \log \frac{\widetilde{M}}{\omega M_v} \right\rceil + i, 0 \right\}$  for some  $i \geq 1$ .

In the first case,  $\mathbb{E}[Y_{a_1^u}Y_{a_2^v}] \leq \frac{\widetilde{M}}{M_v}$ , so the total contribution to the expectation from such vertices is  $\leq \sum_{v \in V(G)} M_v \widetilde{M} \lesssim M^2$ .

In the second case,  $\mathbb{E}[Y_{a_1^u}Y_{a_2^v}] \leq \omega$ , so the total contribution to the expectation from such vertices is  $\leq \sum_{v \in V(G)} M_v^2 \omega = C_1 \omega \leq M^2$ .

In the third case,  $\mathbb{E}[Y_{a_1^u}Y_{a_2^v}] \leq \frac{\widetilde{M}2^i}{M_v}$ , and the probability of  $\mathcal{M}_v$  being at least this high is:

$$\lesssim \min \left\{ 1, \left( \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{kM_v}{2^i \widetilde{M}} \right)^{l-1} \right)^2 \right\} \quad \text{By Lemma 42}$$

Now let  $x \in \mathbb{R}$  be the unique solution to

$$\frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{kM_v}{2^i \widetilde{M}} \right)^{l-1} = 1$$

For  $i \leq x$ ,  $1 \leq 2^{-|i-x|} \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{kM_v}{2^i \widetilde{M}} \right)^{l-1}$ . For  $i \geq x$ ,

$$\left( \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{kM_v}{2^i \widetilde{M}} \right)^{l-1} \right)^2 \leq 2^{-|i-x|} \frac{1}{M_v^2} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{kM_v}{2^i \widetilde{M}} \right)^{l-1}$$

So in either case, the contribution to  $\mathbb{E}[Y^2]$  is bounded by:

$$\begin{aligned}
&\lesssim \frac{\widetilde{M}}{M_v} \sum_{i=\max\{0, -\lceil \log \frac{\widetilde{M}}{\omega M_v} \rceil\}}^{\lceil \log \frac{\sqrt{k}}{\omega} \rceil} 2^{i-|i-x|} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \left( \frac{k M_v}{2^i \widetilde{M}} \right)^{l-1} \\
&\leq \sum_{i=\max\{0, -\lceil \log \frac{\widetilde{M}}{\omega M_v} \rceil\}}^{\lceil \log \frac{\sqrt{k}}{\omega} \rceil} 2^{-|i-x|} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 k^{l-1} \left( \frac{M_v}{2^i \widetilde{M}} \right)^{l-2} \\
&\leq \sum_{i=0}^{\infty} 2^{-i} \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \frac{k^{l-1}}{(2^i \omega)^{l-2}} \\
&\lesssim \sum_{l=2}^s \sum_{S \subseteq V(G), |S|=l, v \in S} M_S^2 \frac{k^{l-1}}{\omega^{l-2}}
\end{aligned}$$

And so, by summing across all  $v$ , the contribution is:

$$\lesssim \sum_{l=2}^s C_l \frac{k^{l-1}}{\omega^{l-2}}$$

Now, by summing across all the cases and using the fact that, with  $\omega \leq \sqrt{k}$  and  $l \geq 2$ ,  $k^{\frac{l}{2}} \leq \frac{k^{l-1}}{\omega^{l-2}}$  (and noting that  $C_l \leq M^2$ ), this gives us:

$$\begin{aligned}
\text{Var}(\overline{M}_H) &\leq \mathbb{E}[Y^2] \\
&\lesssim M^2 + \sum_{l=2}^s C_l \frac{k^{l-1}}{\omega^{l-2}}
\end{aligned}$$

□

### H.3 $\overline{M}$

**Lemma 47.**  $\mathbb{E}[\overline{M}] = M$ .

*Proof.* By Lemma 45,  $\mathbb{E}[\overline{M}|M_L] = \mathbb{E}[\overline{M}_L|M_L] + \mathbb{E}[\overline{M}_H|M_L] = \mathbb{E}[\overline{M}_L|M_L] + \mathbb{E}[\overline{M}_H|M_H] = M_L + M_H = M$ . So  $\mathbb{E}[\overline{M}] = M$ . □

**Lemma 48.**  $\text{Var}(\overline{M}) \lesssim M^2 + \sum_{l=2}^s C_l \left( k^{\frac{l}{2}} + k \left( \frac{C_1^+}{M^2} k \right)^{l-2} \right)$

*Proof.*  $\text{Var}(M) \lesssim \text{Var}(M_L) + \text{Var}(M_H)$ , so from Lemmas 46 and 44, using the fact that  $\frac{k^{l-1}}{\omega^{l-2}} \geq k^{\frac{l}{2}}$  for  $\omega \leq \sqrt{k}$ ,  $l \geq 2$ , we get  $\text{Var}(\overline{M}) \lesssim M^2 + C_2 k \log k + \sum_{l=3}^s C_l \frac{k^{l-1}}{\omega^{l-2}}$ . The result then follows from  $\omega = \min \left\{ \left( \frac{\widetilde{M}}{C_1^+} \right)^2, \sqrt{k} \right\}$ . □

## H.4 Single-Pass Algorithm

**Lemma 49.** *The first and second pass calculations can be performed while storing no more than  $O\left(\frac{m \log \frac{M_V \sqrt{k}}{M^2}}{k}\right)$  edges in expectation.*

*Proof.* We will now demonstrate how both conceptual passes can be calculated in a single pass, by only calculating the value  $\mathcal{M}_v$  for a vertex  $v$  when necessary.

We will store the following edges (noting that edges in this graph are undirected, and so an edge  $e = uv$  is stored if it would be stored as either  $uv$  or  $vu$ ):

$$\begin{aligned} E_0 &= \{uv \in E \mid \exists i, \mathcal{M}_i(u) = 1, \mathcal{M}_i(v) = 1\} \\ E_1 &= \{uv \in E \mid \exists i, c(u) = 1, \mathcal{M}_i(v) = 1\} \\ E_2 &= \{uv \in E \mid c(u) = 1, c(v) = 1\} \\ E_3 &= \{uv \in E \mid \exists i, r_{D,i}(uv) < \frac{2^{h(u)}\omega}{\sqrt{k}}, d(v) = 1\} \\ E_4 &= \{uv \in E \mid r_C(uv) < \frac{2^{h(u)}\omega}{\sqrt{k}}, c(v) = 1\} \end{aligned}$$

Then,  $\forall uv \in E$ :

$$\begin{aligned} \mathbb{P}[uv \in E_0] &\lesssim \frac{1}{k} \\ \mathbb{P}[uv \in E_1] &\lesssim \frac{1}{k} \\ \mathbb{P}[uv \in E_2] &= \frac{1}{k} \\ \mathbb{P}[uv \in E_3] &= \frac{1}{\sqrt{k}} \sum_{i=1}^2 \sum_{h=0}^{\lceil \log \frac{\sqrt{k}}{\omega} \rceil} \mathbb{P}\left[r_{D,i}(uv) < \frac{2^{h(u)}\omega}{\sqrt{k}} \mid h(u) = h\right] \mathbb{P}[h(u) = h] \\ &\lesssim \frac{1}{\sqrt{k}} \sum_{h=0}^{\lceil \log \frac{\sqrt{k}}{\omega} \rceil} \frac{2^h \omega}{\sqrt{k}} \frac{1}{\omega 2^h} \\ &\lesssim \frac{\log \frac{\sqrt{k}}{\omega}}{k} \\ \mathbb{P}[uv \in E_4] &\lesssim \frac{\log \frac{\sqrt{k}}{\omega}}{k} \end{aligned}$$

So these sets will contain  $O\left(\frac{m \log \frac{\sqrt{k}}{\omega}}{k}\right)$  edges in expectation. Then, as  $\omega = \min\left\{\frac{(\widetilde{M})^2}{M_V^\dagger}, \sqrt{k}\right\}$ , and  $M_V^\dagger$  can be any upper bound on  $M_V$ , we can achieve this by storing  $O\left(\frac{m \log \frac{M_V \sqrt{k}}{M^2}}{k}\right)$  edges.

We will then use these to calculate  $\overline{M}_L$  and  $\overline{M}_H$  as follows:

$\overline{M}_L$ : Recall that  $\overline{M}_L$  is the number of copies of  $A$  in  $G'_L$ , the subgraph of  $G$  induced by  $V'_L = \{u \in V(G) | c(u) = 1, \mathcal{M}_u \neq \text{LIGHT}\}$ .  $E_2$  will contain every edge in  $E(G'_L)$ , so we can calculate  $\overline{M}_L$  provided we can calculate  $\mathcal{M}_v$  for all  $v$  s.t.  $c(v) = 1$ .

To do this, we will need to calculate  $\sum_{a \in \mu(\{v\})} X_{v,i}^{(h,a)}$  for  $h = \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil$ , and each  $i \in [2]$ . For each subgraph  $a \in \mu(\{v\})$ ,  $X_{v,i}^{(h,a)} = 1$  iff

$$\begin{aligned} \forall u \in V(a) \setminus \{v\}, \mathcal{M}_i(u) = 1 \\ \forall \{u|vu \in E(a)\}, r_{D,i}(vu) < \frac{\omega 2^h}{\sqrt{k}} \end{aligned}$$

And 0 otherwise. The third and fourth conditions always hold when  $h = \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil$ , so we need to know how many such subgraphs  $a$  exist with  $\forall u \in V(a) \setminus \{v\}, \mathcal{M}_i(u) = 1$ . But in that case,  $\{uw \in E(a) | u \neq v \neq w\} \subseteq E_0$ , and  $\{vu|vu \in E(a)\} \subseteq E_1$ , so  $\sum_{a \in \mu(\{v\})} X_{v,i}^{(h,a)}$  will be equal to the number of subgraphs  $a \in \mu(\{v\})$  s.t.  $\forall u \in V(a) \setminus \{v\}, \mathcal{M}_i(u) = 1$  in the edges we have sampled. So then  $\mathcal{M}_{v,i} = \text{LIGHT}$  iff this number is  $< \frac{\overline{M}_3}{k^{\frac{3}{2}}}$ .

So we can compute  $\mathcal{M}_v$  for each  $v$  s.t.  $c(v) = 1$  using our sampled edges, and therefore we can compute  $\overline{M}_L$ .

$\overline{M}_H$ : For any  $v \in V(G)$ ,  $\overline{M}_v = 0$  if  $h(v) \neq \mathcal{M}_v$ . So it is sufficient to calculate  $\mathcal{M}_v$  when  $h(v) = \mathcal{M}_v$ , and to know that  $h(v) \neq \mathcal{M}_v$  otherwise. We can calculate  $\sum_{a \in \mu(\{v\})} X_{v,i}^{(h,a)}$  for each  $h \leq h(v)$ ,  $i \in [2]$ , as for any  $a \in \mu(\{v\})$  s.t.  $\forall u \in V(a) \setminus \{v\}, \mathcal{M}_i(u) = 1$ ,  $\{uw \in E(a) | u \neq v \neq w\} \subseteq E_0$ , and if  $\forall \{u|vu \in E(a)\}, r_{D,i}(vu) < \frac{\omega 2^h}{\sqrt{k}}$  for some  $h \leq h(v)$ ,  $\{vu|vu \in E(a)\} \subseteq E_3$ .

So then, as  $H_{v,i} = \left\{ h \in \left\{ 1, \dots, \left\lceil \log \frac{\sqrt{k}}{\omega} \right\rceil \right\} \mid X_{v,i}^{(h)} \geq \frac{\omega \overline{M}_2^h}{k^2} \right\}$  we can compute  $H_{v,i} \cap \{0, \dots, h(v)\}$ . As  $\mathcal{M}_{v,i} = \min H_{v,i}$ , we can compute each  $\mathcal{M}_{v,i}$  if it is  $\leq h$ , and determine that it is  $> h(v)$  or LIGHT otherwise. Then, if  $\mathcal{M}_v \leq h(v)$ , at least one of  $\mathcal{M}_{v,1}, \mathcal{M}_{v,2}$  is  $\leq h(v)$ , and so we can calculate it and therefore calculate  $\mathcal{M}_v$ , and if not we can determine that  $\mathcal{M}_v$  is either  $> h(v)$  or LIGHT. (although not necessarily which one)

So if  $\mathcal{M}_v \neq h(v)$ , we know  $\overline{M}_v = 0$ , and then if  $\mathcal{M}_v = h(v)$ , we need to know how many  $a \in \mu(\{v\})$  there are such that:

$$\begin{aligned} \forall u \in V(a) \setminus \{v\}, c(u) = 1 \\ \forall \{u|vu \in E(a)\}, r_C(vu) < \frac{\omega 2^h}{\sqrt{k}} \end{aligned}$$

If these criteria hold, then  $a \subseteq E_2 \cap E_4$ . So we can calculate  $\overline{M}_v$  by counting the number of subgraphs amongst our sampled edges such that these criteria hold.

We then calculate  $\overline{M}_H$  by summing  $\overline{M}_v$  for each  $v \in V(G)$ .

□

We are now ready to prove our generalized theorem.

**Theorem 5** (Subgraph estimation). *Let  $f_\ell$  be the fraction of pairs of subgraphs that intersect at  $\ell$  vertices. We show how to find a  $1 + \epsilon$  factor approximation to  $M$  with probability  $1 - \delta$ , using order*

$$m \frac{\log(1/\delta)}{\epsilon^2} \log M \left( \sum_{\ell=2}^s f_\ell^{2/\ell} + f_\ell^{\frac{1}{\ell-1}} f_1^{1-\frac{1}{\ell-1}} \right)$$

*samples in expectation.*

*Proof.* By Lemma 48, if we choose  $\widetilde{M} \leq M$ , we can obtain variance  $\lesssim M^2$  with

$$k = \sum_{l=2}^s \left( \left( \frac{M^2}{C_l} \right)^{\frac{2}{l}} + \left( \frac{M^2}{C_l} \right)^{\frac{1}{l-1}} \left( \frac{M^2}{C_1} \right)^{\frac{l-2}{l-1}} \right) = \sum_{l=2}^s \left( \left( \frac{M^2}{C_l} \right)^{\frac{2}{l}} + \frac{M^2}{(C_l C_1^{l-2})^{\frac{1}{l-1}}} \right)$$

, by taking our upper bound  $C_1^+$  within a constant factor of the true value of  $C_1$ .

So as by Lemma 49, we need to sample  $\frac{\log k}{\epsilon}$  edges to attain this, and we can then repeat the algorithm  $O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta}\right)$  times, taking a median-of-means of our results, to get an estimate within  $\epsilon M$  of  $M$  with probability  $1 - \delta$ .

Our result then follows from the fact that  $f_l = \Theta\left(\frac{C_l}{M^2}\right)$ . □

**Corollary 50.** *We can obtain a constant-error approximation to  $M$  while keeping*

$$O\left(m \sum_{l=2}^s \left( f_l^{\frac{2}{l}} + f_l^{\frac{1}{l-1}} f_1^{1-\frac{1}{l-1}} \right)\right)$$

*edges.*

## I Transitivity coefficient

Let  $P_2$  denote the number of length 2 paths in the graph, so the transitivity coefficient  $\alpha = 3T/P_2$ . The condition that the graph have no isolated edges implies that  $P_2 \gtrsim m$ , or  $\alpha \lesssim T/m$ . Hence to show that the bound in [BFKP14] is weaker than the one in [PT12], it suffices to show the following lemma.

**Lemma 51.** *Consider any graph with  $m$  edges,  $T$  triangles,  $P_2$  length-2 paths, transitivity coefficient  $\alpha = 3T/P_2$ , and at most  $\Delta_E$  triangles sharing a common edge. Then*

$$\frac{\sqrt{m}}{\alpha} + \frac{m\sqrt{m}}{T} \gtrsim m \left( \frac{\Delta_E}{T} + \frac{1}{\sqrt{T}} \right)$$

*Proof.* First, consider the  $\frac{m\Delta_E}{T}$  term on the right. This is trivially bounded by the left if  $\Delta_E \leq \sqrt{m}$ . Otherwise, note that because  $\Delta_E$  triangles share a common edge, the corresponding vertices have degree at least  $\Delta_E + 1$  so  $P_2 \geq 2\binom{\Delta_E+1}{2} \geq \Delta_E^2$ . Hence

$$\frac{\sqrt{m}}{\alpha} = \frac{\sqrt{m}P_2}{3T} \geq \frac{\sqrt{m}\Delta_E^2}{3T} \geq \frac{m\Delta_E}{3T}$$

from  $\Delta_E > \sqrt{m}$ , so  $\frac{m\Delta_E}{T}$  is bounded by the LHS.

Now, consider bounding the  $m/\sqrt{T}$  term. If  $T \geq m$ , then  $P_2 \geq 3T \geq 3m$ , so

$$\frac{\sqrt{m}}{\alpha} = \frac{\sqrt{m}P_2}{3T} \geq \frac{m}{\sqrt{T}} \sqrt{\frac{P_2}{3T}} \geq \frac{m}{\sqrt{T}}.$$

On the other hand, for  $T \leq m$ ,  $m/\sqrt{T} \leq \frac{m\sqrt{m}}{T}$ .

Hence both terms of the RHS are bounded by the LHS, so the sum is bounded by twice the LHS.  $\square$

The result of [JSP13] for triangle counting is also implied by the [PT12] bound. In [JSP13],  $\frac{m}{\epsilon^2\sqrt{T}}$  space is used to learn  $\alpha$  to  $\pm\epsilon$ ; this gives a sample complexity of

$$\frac{m}{\sqrt{T}} \cdot \left(\frac{P_2}{T}\right)^2$$

for learning a constant multiplicative approximation to  $T$ .

**Lemma 52.** *Consider any graph with  $m$  edges,  $T$  triangles,  $P_2$  length-2 paths, transitivity coefficient  $\alpha = 3T/P_2$ , and at most  $\Delta_E$  triangles sharing a common edge. Then*

$$\frac{m}{\sqrt{T}} \cdot \left(\frac{P_2}{T}\right)^2 \gtrsim m \left(\frac{\Delta_E}{T} + \frac{1}{\sqrt{T}}\right)$$

*Proof.* The  $\frac{m}{\sqrt{T}}$  term on the right is trivial, since  $T \lesssim P_2$ . For the other term, similarly to the previous proof, we have

$$\frac{m}{\sqrt{T}} \cdot \left(\frac{P_2}{T}\right)^2 \gtrsim \frac{m}{\sqrt{T}} \cdot \left(\frac{P_2}{T}\right)^{.5} = \frac{m\sqrt{P_2}}{T} \geq \frac{m\Delta_E}{T}.$$

$\square$