

# Explicit, Closed-form, General bounds for Cuckoo Hashing with a Stash

Daniel Noble

University of Pennsylvania, dgnoble@cis.upenn.edu

November 17, 2022

## Abstract

Cuckoo Hashing is a dictionary data structure in which a data item is stored in a small constant number of possible locations. It has the appealing property that a data structure of size  $2m$  can hold up to  $n = \frac{1}{d}m$  elements for any constant  $d > 1$ ; *i.e.*, the data structure size is a small constant times larger than the combined size of all inserted data elements. However, the probability that a cuckoo hash table build fails is  $\Theta(\frac{1}{m})$ . This is too high for many applications, especially cryptographic applications and Oblivious RAM. An alternative proposal introduced by Kirsch et al. is to store elements which cannot be placed in the main table in a “stash”, reducing the failure probability to  $\mathcal{O}(m^{-(s+1)})$  where  $s$  is any constant stash size. However, this analysis did not apply to super-constant  $s$ , and the bounds are asymptotic rather than concrete. Further works improved upon this, but either were not concrete, not closed-form or had limitations on the stash size. In this paper we present the first concrete, closed-form bounds for the failure probability of cuckoo hashing with a stash for general stash sizes.

## 1 Overview

This paper proves the first explicit, closed-form bounds for the failure probability of cuckoo hashing with a stash for general stash sizes. Specifically, it proves the following bound:

**Theorem 1.** *Given a 2-table cuckoo hash table, in which each table holds  $m$  elements, and which has a stash of size  $s \geq 0$ , the cuckoo hash table can successfully hold  $n \leq \frac{1}{d}m$  elements for any  $d > 1$ , with failure probability at most:*

$$C(s+2) \left( \frac{c(s+1)}{m} \right)^{s+1}$$

where  $C$  and  $c$  depend only on  $d$ :

$$C = \frac{16e^2 d^3 (d+1)^2}{(d-1)^5} e^{\frac{7.2d^3(d+1)^2}{(d-1)^5}}$$

$$c = \frac{d^2(d+1)^2}{e(d-1)^5}$$

As a corollary, this shows tighter bounds than what existed previously on the requirements on  $m$  and  $s$  to have failure probability negligible in some term  $N$ . Specifically:

**Corollary 1.** *The probability of build failure for a 2-table cuckoo hash table of size  $m = \omega(\log(N))$  and a stash of size  $s = \Theta(\log(N))$  is negligible in  $N$ .*

The analysis takes the standard approach of representing the problem as a problem on random bipartite graphs. The result above is attained by proving the first explicit closed-form general upper bound for the excess of balanced bipartite graphs with randomly chosen edges. Namely, it proves the following:

**Theorem 2.** *Let  $G$  be a bipartite multigraph with parts  $L$  and  $R$ , where  $|L| = |R| = m$ , and  $n = \frac{1}{d}m$  edges each chosen independently and uniformly at random from  $L \times R$ . Let  $\mathbf{ex}(G)$  represent the excess of  $G$ , that is the minimum number of edges that need be removed such that every connected component has at most one cycle. Then:*

$$\Pr(\mathbf{ex}(G) \geq t) \leq C(t+1) \left(\frac{ct}{m}\right)^t$$

for every  $t \geq 1$  and  $d > 1$  where  $C$  and  $c$  are defined as above.

The paper is organized as follows. Section 2 presents a summary of previous work. Section 3 explains cuckoo hashing and Section 4 explains the correspondence between cuckoo hashing and random bipartite graphs. Section 5 presents the main analysis on the graph problem, though proofs for some of the lemmas for this analysis are deferred to Section 6. As a bonus, Section 7 shows that Corollary 1 is basically tight; that is if  $m = \mathcal{O}(\log(N))$  it is essentially impossible to attain negligible failure probability in  $N$ .

## 2 Previous Work

Cuckoo hashing is a hash table implementation that improves performance by allowing objects to be stored in a number of locations [PR01]. Pagh and Rodler discovered that this small modification greatly reduced the probability of a build failure. Specifically, a cuckoo hash table can store  $n$  elements of size  $W$  in  $\Theta(nW)$  space, with accesses only accessing  $\Theta(1)$  locations of size  $W$ . It is still sometimes impossible to place all items in the required locations, in which case the build fails. A build failure, which we henceforth refer to just as failure, occurs with probability  $\Theta(\frac{1}{n})$  (see [DK12] for the explicit constant).

For many applications this failure probability is sufficient. In the case that a failure does occur, the hash table can simply be rebuilt with new hash functions. While a hash table rebuild requires  $\Theta(n)$  computation, the probability

of this occurrence is  $\Theta(\frac{1}{n})$ , so the *amortized* computation cost per access is still constant.

However, for many applications this failure probability is still too high. In particular, in security-oriented applications, rebuilds often constitute a security failure. A sequence of works arose to reduce the failure probability of Cuckoo hashing.<sup>1</sup>

First, Kirsch, Mitzenmacher and Wieder introduced the modification that any items which could not be stored in the Cuckoo Hash table would be stored in a “stash” of constant size  $s$  [KMW09]. They showed that the probability of a build failure<sup>2</sup> was then reduced to  $\mathcal{O}(n^{-(s+1)})$ . (They state this as the equivalent statement that the probability that the required stash is at least  $s$  is  $\mathcal{O}(n^{-s})$ .) This analysis allowed the failure probability to be reduced significantly, but it only applied to constant  $s$ . Furthermore, closed-form explicit bounds are not stated.

Kutzelnigg provided a detailed analysis of cuckoo hashing with a stash based on a generating function approach [Kut10]. In some senses the approach is very theoretically satisfying. The paper shows (Theorem 1) that the failure probability is exactly

$$c(\alpha, s)m^{-s-1} - \mathcal{O}(m^{-s-2})$$

where  $m$  is the size of each table,  $s$  is the stash size,  $\alpha$  is related to the ratio between the size of the table and the number of entries and  $c$  depends only on  $\alpha$  and  $s$ , where  $c(\alpha, s) \neq 0$ . However, while the paper states that  $c(\alpha, s)$  can be calculated explicitly using the techniques in the paper, this calculation is very complex. In particular, it states: “The calculations [to compute  $c(\alpha, s)$ ] are limited by the available memory of the machine that executed the computer algebra system. Using a workstation with 12GB RAM, we were successful in solving the problem for  $s \in \{0, 1, 2\}$ .” Thus, it seems that a standard workstation is not able to compute  $c(\alpha, s)$  for  $s > 2$ . Thus, while the paper presents powerful techniques and proves that some constant  $c(\alpha, s)$  exists, it cannot yet be considered a practical technique for calculating good upper bounds on the failure probability for general values of  $s$ .

Goodrich and Mitzenmacher developed another analysis for super-constant values of  $s$  in order to generate an improved Oblivious RAM protocol [GM11]. Oblivious RAM is a technique to hide virtual accesses from an adversary who can see physical accesses, and numerous protocols have used cuckoo hashing for this purpose (*e.g.*, [PR10], [GMOT12], [KLO12]). In these protocols, cuckoo hashing build failures constitute a security failure. Failure is desired to be negligible in some parameter  $N$ , but tables can be much smaller than  $N$ , for instance they may be polylogarithmic in  $N$ . Goodrich and Mitzenmacher extended the

<sup>1</sup>Since this paper examines the 2-table version, (where there are 2 tables, each of size  $m = dn$ ) we only discuss previous work with results for this version. However, as far as the authors are aware, there are also no explicit, closed-form bounds for general stash sizes for other versions of cuckoo hashing.

<sup>2</sup>While this paper presents failure probabilities in terms of  $m$ , if other papers present their results in terms of  $n$ , we do too. Since  $n = \frac{1}{d}m$  converting between these is trivial.

analysis of Kirsch et al. to achieve negligible failure probability in certain cases. They proved that, provided  $n = \Omega(\log^7(N))$ , the probability of a build failure is upper bounded by  $n^{-\Omega(s)}$  for general values of  $s$ . Thus, a stash of size  $s = \Theta(\log(N))$  would result in a failure probability negligible in  $N$ . In comparison, our work does not have limits on the value of  $n$ , and states explicit bounds.

Aumüller, Dietzfelbinger and Woelfel then presented an elegant alternative analysis of cuckoo hashing with a stash based on graph counting [ADW14]. In this analysis, they showed a new result for super-constant stash sizes. They showed that for sufficiently large  $n$ , the failure probability is  $\mathcal{O}(n^{-\frac{\varrho}{2}})$  when  $s \leq n^{\frac{1}{3e}}$ , for a suitable constant  $\varrho$ . However this constraint on the relationship between  $s$  and  $n$  proves restrictive in practice. They state “a rough estimate ... shows that  $\varrho = 27$  suffices”, but this would impose the restriction that  $s \leq n^{\frac{1}{81}}$ . Even for  $s = 2$ , this makes the bound inapplicable for practical values of  $m$ .

Pinkas et al. [PSSZ15] show empirically that cuckoo hashing with a stash performs well even for small values of  $d$ . Specifically, they generated random instances of the problem for  $n \in \{2^{11}, 2^{12}, 2^{13}, 2^{14}\}$  and  $d = 1.2$ . For each value of  $n$  they generated  $2^{30}$  random instances of the cuckoo hashing table and determined the stash sizes that were needed. From this they extrapolated the stash sizes needed for smaller error probabilities and other values of  $n$  (see [PSSZ15] Table 3 and Figure 1). Compared to their work, our result has the following advantages. Firstly, it is based on mathematical proofs rather than empirical analysis. Secondly, the results of Pinkas et al. for small error probabilities (e.g.  $2^{-40}$ ) depend on extrapolations.<sup>3</sup> Thirdly, Pinkas et al. only show the failure probabilities for specific values of  $n$ ,  $d$  and  $s$ , whereas the results of this paper show explicit upper bounds for the failure probability for any values of  $n$ ,  $d$  and  $s$ . On the negative side, compared to Pinkas et al. our constants are clearly worse. Below we present parameter choices that, based on our analysis, have failure probability  $2^{-40}$  when  $d = 3$ . Compare these to Table 4 of [PSSZ15]. When  $d = 3$ , our memory usage is 2.5x that of [PSSZ15] and the required stash size is about 4x larger.

<b>number of elements</b> $n$	$2^{10}$	$2^{12}$	$2^{14}$	$2^{16}$	$2^{18}$	$2^{20}$	$2^{22}$	$2^{24}$
<b>stash size</b> $s$	33	23	18	15	12	11	9	8

Table 1: Stash sizes needed to obtain failure probability below  $2^{-40}$  when  $d = 3$ .

Finally, Pinkas et al. [PSWW18a] observed that an analysis by Wieder [Wie16] provided a bound on the stash size that does not make any assumptions about the stash size (see Appendix C of the full version [PSWW18b]). While Wieder’s proof does not present explicit bounds in the final result, Pinkas et al. filled in the missing details to present a explicit, general bound. However,

<sup>3</sup>This limitation is, in a sense, inherent in empirical evidence of statistical failure. We would like, ideally, to pick some upper bound on the number of times the program will be run and argue that it is unlikely to fail even if run that many times. But to prove this, we need to run it that many times in the lab, and it is usually unrealistic to expect the program to be run more times in a lab than in production.

the bound they present is still not closed-form and it also seems that Wieder’s proof contained an error which propagated to the bound of Pinkas et al.<sup>4</sup>

We could make the bound closed-form, such as by applying Lemma 7 from this paper. This would result in the following bound for the failure probability:

$$\left(\frac{2}{m}\right)^{s+1} 2e \left(\frac{10(s+1)+1}{e \ln(d)}\right)^{10(s+1)+1}$$

Asymptotically, for non-constant  $s$ , the bound presented in Theorem 1 of this paper is much tighter. The bound above can be represented asymptotically as:

$$\mathcal{O}(\mathcal{O}(1)^s m^{-(s+1)} s^{10s+11})$$

whereas the equivalent representation of Theorem 1 would be:

$$\mathcal{O}(\mathcal{O}(1)^s m^{-(s+1)} s^{s+2})$$

### 3 Cuckoo Hashing

Cuckoo Hashing in its simplest form involves 2 hash functions,  $h_1$  and  $h_2$ , and 2 hash tables,  $T_1$  and  $T_2$ , each with  $m = dn$  locations of capacity 1. Each hash table has a unique hash function, and the hash functions are assumed to produce outputs uniformly at random in  $\{1, \dots, m\}$ . The tables consist of pairs  $(x, y)$  where  $x$  is the dictionary key and  $y$  is the dictionary value. An item  $(x, y)$  is stored in the table by being inserted into  $T_1[h_1(x)]$ . If another item  $(x', y')$  was stored in that location, it is removed from its original location (like a baby bird being displaced from its nest by a Cuckoo chick) and is placed in  $T_2[h_2(x')]$ . This may replace another item, which the algorithm likewise attempts to insert. This process continues either until every item has found a location in which to be inserted or it is determined that it is impossible to place all items in the cuckoo hash table.<sup>5</sup> In the latter case the insertion has “failed”. This triggers

<sup>4</sup>Wieder’s proof uses the inductive hypothesis that the number of non-isomorphic graphs with  $t$  edges,  $k$  components and excess  $s$  is at most  $t^{5s+2k}$ , but based on the previous parts of the proof this should be  $t^{8s+2k}$ , which, since  $k \leq s$ , is at most  $t^{10s}$ , not  $t^{8s}$ . This means that the explicit, non-closed-form bound that the excess is at least  $s$ , that would be obtained from Wieder’s proof should be:

$$\frac{1}{n^s} \left(\frac{2}{1+\epsilon}\right)^s \cdot \sum_{t \leq n} \frac{t^{10s}}{(1+\epsilon)^t} = \left(\frac{2}{m}\right)^s \cdot \sum_{t \leq n} \frac{t^{10s}}{d^t}$$

where  $d = (1 + \epsilon) = \frac{m}{n}$ .

<sup>5</sup>Some works fix a maximum recursion depth for the insertion procedure, such as [KMW09] which sets it to  $\alpha \log(N)$  for a sufficiently large constant  $\alpha$ . It is then possible, with a small but non-negligible (in  $N$ ) probability, that an insertable item is not inserted. Our analysis instead assumes an optimal allocation. This can be achieved by only stopping the insertion when it is detected that the insertion process has entered an infinite loop. See Section 4 of [Kut10] for more details. This process still has an expected insertion time of  $\Theta(1)$  per item (Lemma 3.7 of [Aum10]), so building the full table this way takes expected  $\Theta(n)$  time.

a “table rebuild” in which new tables are created with new hash functions and the algorithm attempts to insert every element into the new hash table.

In cuckoo hashing with a stash, if an item cannot be inserted, the build does not immediately fail, but instead the item is placed in the stash. The build only fails if an item cannot be inserted and the stash is already full.

Cuckoo hash tables can be generalized to have a larger constant number of tables, or have locations with some constant capacity greater than 1. They can also be generalized to use multiple hash functions in a single table. However, this work will analyze only the traditional 2-table version.

## 4 Graph Representation

Analyses of cuckoo hash table failure often represent the problem as a graph problem as follows. For each location in the cuckoo hash table, create a vertex. Since the cuckoo hash table has two tables each of size  $m$ , there will be  $2m$  vertices. For each element stored in the cuckoo hash table, draw an edge between the two locations in which it may be stored, so  $n$  edges total. Let  $G$  be the resulting graph. Since there will be one location from each table,  $G$  will be bipartite, with  $m$  vertices in each part. There may also be multiple edges between a pair of vertices, so  $G$  is a multigraph. Observe also that the graph is not connected: since  $n < m$  some nodes will not be connected to any edges and there may also be multiple connected components that contain edges.

We introduce some graph notation and terminology. Given a graph  $G$ , let  $\gamma(G)$  denote the cyclotomic number of  $G$ , that is the minimum number of edges that must be removed in order for  $G$  to have no cycles. Let  $\mathbf{ex}(G)$  denote the *excess* of  $G$ , that is the minimum number of edges that must be removed from  $G$  to ensure that every connected component has at most one cycle.

Analysis is based on the following critical observation (which is proven, for instance, as Lemma 5 of [ADW14]) that build failure in a cuckoo hash table can be determined by the excess of its graph representation:

**Theorem 3.** *Let  $G$  be the graph representation of a cuckoo hash table with a stash of size  $s$ , where  $s$  is any non-negative integer. Then the build succeeds if and only if  $\mathbf{ex}(G) \leq s$ . Equivalently, the build fails if and only if  $\mathbf{ex}(G) \geq s + 1$ .*

Now, let  $G(m, m, n)$  be the distribution of graphs generated from graph representations of random cuckoo hash tables with 2 tables of size  $m$  and  $n$  items. Let  $G \leftarrow G(m, m, n)$ , i.e.,  $G$  is randomly sampled from  $G(m, m, n)$ . As already stated,  $G$  will then be a bipartite graph with parts  $A$  and  $B$  each of  $m$  vertices. Since each hash function produces a random value in  $\{1, \dots, m\}$ ,  $G$  will have  $n$  edges chosen uniformly at random from  $A \times B$ . Note that this exactly matches the description of how  $G$  is chosen in Theorem 2. Therefore, Theorem 2 describes the bounds on the excess of graph representations of random cuckoo hash tables. Thus, by Theorem 3, Theorem 2 implies Theorem 1.

## 5 An Explicit Analysis

The remainder of this section proves Theorem 2, (though some more tedious lemmas are deferred to Section 6).

The analysis will proceed by observing the randomly generated graph, carefully choosing which information about the graph is revealed. This will give bounds on both the size and the cyclotomic number of the component containing a randomly chosen edge, which will give bounds on the excess of the entire graph. We will make extensive use of the following algorithm, which should be viewed as occurring on a randomly generated graph. Observed variables are therefore drawn from probability distributions and the distribution of the remainder of the graph at any point is conditioned on the variables that have already been observed.

### Edge Component Search Algorithm

1. While  $\exists$  an undiscovered edge in  $G$ 
  - (a) Select one such edge at random. Call it  $e$ . Call the vertices it connects  $v_1$  and  $v_2$ .
  - (b) Let  $Q_e$  be a queue initialized to  $\{v_1, v_2\}$ . Set  $V_e = \{\}$ ,  $Y_e = \{e\}$ .
  - (c) While  $Q_e$  is not empty
    - i. Set  $v \leftarrow \text{dequeue}(Q_e)$ .
    - ii. Add  $v$  to  $V_e$ .
    - iii. Set  $N_v$  to be the set of undiscovered neighbors of  $v$ , (*i.e.*, neighbors of  $v$  which have never been placed in  $Q_e$  for any edge  $e$ ). This should be thought of as first observing  $|N_v|$ , and then observing the vertices themselves.
    - iv. Enqueue all vertices in  $N_v$  to  $Q_e$ .
    - v. For each  $w$  in  $N_v$ , add one of the edges connecting  $v$  to  $w$  to  $Y_e$ . (If there is more than one such edge, pick one at random, without observing the total number of such edges.)
  - (d) Set  $T_e = (V_e, Y_e)$ .
  - (e) For every pair of vertices in  $V_e$  which are in different parts of  $G$ , observe the number of unobserved edges between the vertices. Set  $Z_e$  to be the set of these edges.
  - (f) Set  $C_e = (V_e, Y_e \cup Z_e)$ .

**Theorem 4.**  $C_e$  calculated in step 1f will be the connected component in  $G$  containing the edge  $e$  chosen in step 1a.

*Proof.* Observe that steps 1b and 1c are identical to a Breadth First Search (BFS), except that the queue begins containing two vertices instead of 1. However,  $v_2$  is a neighbor of  $v_1$ , so the initial state of the system can be viewed as the state of a BFS starting at  $v_1$  where the first neighbor ( $v_2$ ) has already been

found, and added to the queue. Therefore, the resulting BFS will find exactly the nodes reachable from  $v_1$  in  $G$ , which is exactly the nodes in the connected component in  $G$  containing  $e$ . Observe also that any edges that may exist in  $C_e$  are found, either in steps 1b and 1(c)v, in which case they are added to  $Y_e$ , or in step 1e in which case they are added to  $Z_e$ . Either way, these edges exist in  $C_e$ . Lastly, only edges in the connected component containing  $e$  exist in  $C_e$ , since only edges in  $G$  connecting vertices in  $C_e$  are added.  $\square$

Observing the correspondence to a BFS also indicates the following Lemma.

**Lemma 1.**  $T_e$  calculated in step 1d is a spanning tree of  $C_e$ .

Furthermore,  $Y_e$  and  $Z_e$  are disjoint. Therefore,  $Z_e$  contains a set of edges in  $C_e$  which, if removed, produces a tree. This implies the following fact.

**Fact 1.**  $|Z_e| = \gamma(C_e)$

Lastly, since the random edge selected in 1a will always be one that has not yet been discovered, and all edges in a component are discovered when that component is explored, each new component found will be separate to all previous components found. Furthermore, since the algorithm continues until all edges are found, it will find all components of  $G$ .

We will now prove some claims about the distributions of neighbors, and edges, found by the Edge Component Search Algorithm. The analysis is made challenging by the fact that the existence of edges is not independent: the more edges are found, the fewer are left to find. Likewise, the more vertices are explored by the BFS, the fewer are left as possible end-points for the remaining edges. However this is resolved, in short, due to the way variables are observed in the Edge Component Search Algorithm, as edges are discovered at least as quickly as vertices are found in either part. Thus, even though the actual probability distributions depend on what has already been discovered, we can find probability distributions that only depend on  $m$  and  $n$  that stochastically dominate the real distributions.

We first upper-bound the number of neighbors found in step 1(c)iii.

**Theorem 5.** *The number of neighbors found in step 1(c)iii is stochastically dominated by  $\text{Bin}(n, \frac{1}{m})$*

*Proof.* First, the number of undiscovered neighbors of  $v$  found in step 1(c)iii is at most the number of undiscovered edges that connect to  $v$ . Let  $u$  be the number of discovered edges at a certain point of time, and  $n - u$  be the number of undiscovered edges. Let  $A$  be the part of the bipartite graph containing  $v$  and  $B$  the other part. Each of the  $u$  edges has one end-point in  $A$ . One of these  $u$  edges is known to have its end-point in  $A$  at  $v$  (for  $v_1, v_2$  this is  $e$ , and for other vertices, it is the edge that was used to find  $v$ ). Therefore, there are  $u - 1$  edges that have other end-points in  $A$ , and so at most  $u - 1$  vertices in  $A$  that are end-points of previously discovered edges. Only vertices that are end-points of a previously-discovered edge may have their number of neighbors examined (in



step 1(c)iii). Therefore, there are at least  $m - u + 1$  vertices in  $A$  (including  $v$ ) which prior to step 1(c)iii have not had their number of neighbors examined.

Some of the remaining  $n - u$  edges may be later discovered to exist between previously found vertex-pairs (in step 1e). The number of undiscovered edges that are not in this category is still at most  $n - u$ .

Therefore, there are at most  $n - u$  edges that could contribute towards  $|N_v|$ , and for each, the only thing that is known about the edge's end-point in  $A$  is that it is not one of the at most  $u - 1$  vertices in  $A$  which have had their number of neighbors counted. Hence, each such edge will have  $v$  as its end-point in  $A$  with probability at most  $\frac{1}{m-u+1}$ . Since there are at most  $n - u$  such edges,  $|N_v|$  is stochastically dominated by  $\text{Bin}(n - u, \frac{1}{m-u+1})$ , which by Lemma 3 is stochastically dominated by  $\text{Bin}(n, \frac{1}{m})$  □

We now need to show bounds on the number of edges found in step 1e. It will help to first define three types of vertex pairs. The first are opened vertex pairs, for which the number of edges between the pair of vertices is fully known (including when it is known to be zero). Step 1e cannot find any edges between opened vertex pairs, since it only finds previously undiscovered edges. The second type is partially opened vertex pairs, for which it is known that at least one edge exists between them but it is not known how many more exist. The third type is unopened vertex pairs, for which it is not yet known whether the vertices are neighbors.

We begin by showing bounds on the number of edges between unopened vertex pairs.

**Theorem 6.** *In the Edge Component Search Algorithm, if at a point in time  $u$  edges have been discovered, then the number of edges between an unopened vertex pair,  $v$  and  $w$ , is stochastically dominated by  $\text{Bin}(n, \frac{1}{m(m-u)})$ .*

*Proof.* First we show that the number of edges between  $v$  and  $w$  is stochastically dominated by  $\text{Bin}(n - u, \frac{1}{(m-u)^2})$ .

Let  $q$  be the number of edges that exist between unopened vertex pairs. Every such edge must not yet have been discovered, but there may be some undiscovered edges between partially opened vertex pairs. Therefore  $q \leq n - u$

For the  $q$  edges that exist between unopened vertex pairs, we do not know any information about which vertices they exist between beyond the fact that they exist between unopened vertex pairs. Furthermore, it is equally likely to exist between any such pair.

Since only  $u$  edges have been discovered, there must be at least  $m - u$  vertices in each part that touch no discovered edges. Hence, each pair of such vertices is an unopened vertex pair. Therefore, there are at least  $(m - u)^2$  unopened vertex pairs. Thus, for any given unopened pair, and an edge that exists between a unopened pair, the probability that the edge exists between that unopened pair is at most  $\left(\frac{1}{m-u}\right)^2$ . Hence, the number of edges between any given unopened

pair will be stochastically dominated by  $\text{Bin}(n - u, \frac{1}{(m-u)^2})$ , which by Lemma 4, is stochastically dominated by  $\text{Bin}(n, \frac{1}{m(m-u)})$ .  $\square$

Next we show bounds on the number of additional edges found in step 1e between partially opened pairs.

**Theorem 7.** *In the Edge Component Search Algorithm, the number of additional edges between a partially opened vertex pair is stochastically dominated by the number of edges between an unopened vertex pair.*

(Proof deferred to section 6).

Combining this with Theorem 6 and observing that  $u \leq n$ , we get the following result.

**Theorem 8.** *In the Edge Component Search Algorithm, any vertex pair that is partially opened or unopened, has a number of undiscovered edges that is stochastically dominated by  $\text{Bin}(n, \frac{1}{m(m-n)})$ .*

Now define a function  $H(m, n)$ , which samples a graph  $H \leftarrow H(m, n)$  chosen the same as  $C_e$  in step 1 of the Edge Component Search Algorithm except that:

- Edges and vertices are given new unique identifiers when discovered that may not be the same as the names “found” by the Edge Component Search Algorithm.
- $|N_v|$  in step 1(c)iii is chosen from  $\text{Bin}(n, \frac{1}{m})$
- In step 1e, the additional edges between any pair of vertices in different parts is chosen from  $\text{Bin}(n, \frac{1}{m(m-n)})$ . (Recall the graph is a tree at this point, so is bipartite.)
- We refer to  $V_e$  as  $V$ ,  $Q_e$  as  $Q$ ,  $Y_e$  as  $Y$  and  $Z_e$  as  $Z$ .

**Theorem 9.** *For any component  $C_e$  discovered in the Edge Component Search Algorithm,  $\gamma(C_e)$  is stochastically dominated by  $\gamma(H)$  for an independent sample  $H \leftarrow H(m, n)$ .*

*Proof.* We can view the two graph-sampling algorithms as running in parallel using the same source of randomness. We can choose an interpretation of the randomness generated such that if an event in the sampling of  $H$  stochastically dominates an event in the sampling of  $C_e$ , the event always happens in  $H$  if it happens in  $C_e$ . Since the probability of finding an edge in  $C_e$  is always stochastically dominated by that of finding the edge in  $H$  (from Theorems 5 and 8),  $C_e$  will be a subset of  $H$  for any choice of randomness. Therefore  $\gamma(C_e) \leq \gamma(H)$  for any choice of randomness, which implies that  $\gamma(C_e)$  is stochastically dominated by  $\gamma(H)$ .  $\square$

We can now upper bound  $|H|$  and  $\gamma(H)$  in order to upper bound  $\gamma(C_e)$ .

**Theorem 10.** For  $H \leftarrow H(m, n)$ , where  $m = dn$  for  $d > 1$ , and  $c_3 = \frac{(d-1)^2}{d+1}$ , for  $k \geq 2$ ,

$$\Pr(|H| \geq k) \leq \frac{2d^2}{k-1} e^{-c_3 k}$$

*Proof.* Now, the vertices of  $H$  are found by each vertex having a number of children chosen from the distribution  $\text{Bin}(n, \frac{1}{m})$ . Therefore  $(V, Y)$  can be viewed as the result of a Galton-Watson Branching process, with 2 roots, and children chosen from independent samples of  $\text{Bin}(n, \frac{1}{m})$ . The Otter-Dwass formula [Pit98, Dwa69] states that the probability that a Galton-Watson process that initially has  $\alpha$  nodes, will be of size  $k$  is exactly

$$\frac{\alpha}{k} \Pr(S_k = k - \alpha)$$

where  $S_k$  is the distribution of  $k$  samples of the progeny distribution. In this case  $S_k = \text{Bin}(nk, \frac{1}{m})$ . Therefore, for  $k \geq 2$

$$\begin{aligned} \Pr(|H| = k) &= \frac{2}{k} \Pr(\text{Bin}(nk, \frac{1}{m}) = k - 2) \\ &= \frac{2(k-1)}{(kn-k+2)(kn-k+1)} m^2 \left(1 - \frac{1}{m}\right)^2 \Pr(\text{Bin}(nk, \frac{1}{m}) = k) \\ &\leq \frac{2d^2(k-1)}{(k - \frac{k}{n})(k - \frac{k}{n})} \Pr(\text{Bin}(nk, \frac{1}{m}) = k) \\ &\leq \frac{2d^2}{k-1} \Pr(\text{Bin}(nk, \frac{1}{m}) = k) \\ &\leq \frac{2d^2}{k-1} \Pr(\text{Bin}(nk, \frac{1}{m}) \geq k) \\ &\leq \frac{2d^2}{k-1} e^{-c_3 k} \end{aligned}$$

where  $c_3 = \frac{(d-1)^2}{d(d+1)}$ . The last step comes from the Chernoff Bound

$$\Pr(X \geq (1 + \delta)\mu) \leq e^{-\delta^2 \mu / (2 + \delta)}$$

where the expected value  $\mu = \frac{k}{d}$ . □

Now we can bound  $\gamma(H)$  for a given  $|H|$ .

**Theorem 11.** For  $H \leftarrow H(m, n)$ ,

$$\Pr(\gamma(H) \geq t | |H| = k) \leq \left( \frac{enk^2}{4m(m-n)t} \right)^t$$

*Proof.*  $H$  is bipartite. If one part has size  $a$ , the other has size  $k - a$ . The cyclotomic number of  $H$  is the number of additional edges added in the last step. The number of pairs of vertices that may have edges added between them is  $a(k - a)$  which has maximum value  $\lfloor \frac{k}{2} \rfloor \lceil \frac{k}{2} \rceil = \lfloor \frac{k^2}{4} \rfloor$ .

Each such vertex pair has a number of edges drawn from the distribution  $\text{Bin}(n, \frac{1}{(m-n)m})$ . Therefore the total number of edges is stochastically dominated by  $\text{Bin}(n \lfloor \frac{k^2}{4} \rfloor, \frac{1}{m(m-n)})$ .

Applying the Chernoff bound from Lemma 6 completes the proof.  $\square$

**Theorem 12.** *For any component  $C_e$  found by the Edge Component Search Algorithm,*

$$\Pr(\gamma(C_e) \geq t) \leq c_4 \left( \frac{c_5 t}{m} \right)^t$$

where  $c_4 = 8ed^2$  and  $c_5 = \frac{1}{e(d-1)c_3^2} = \frac{d^2(d+1)^2}{e(d-1)^5}$ .

*Proof.* First we show bounds on  $\gamma(H)$ , where  $H \leftarrow H(m, n)$ . Combining Theorem 10 and Theorem 11 we can obtain bounds for  $\gamma(H)$ :

$$\begin{aligned} \Pr(\gamma(H) \geq t) &\leq \sum_{k=2}^{\infty} \Pr(|H| = k) \Pr(\gamma(H) \geq t | |H| = k) \\ &\leq \sum_{k=2}^{\infty} \frac{2d^2}{k-1} e^{-c_3 k} \left( \frac{enk^2}{4m(m-n)t} \right)^t \\ &\leq 2d^2 \left( \frac{en}{4m(m-n)t} \right)^t \sum_{k=2}^{\infty} \frac{1}{k-1} e^{-c_3 k} k^{2t} \\ &\leq 4d^2 \left( \frac{en}{4m(m-n)t} \right)^t \sum_{k=2}^{\infty} e^{-c_3 k} k^{2t-1} \end{aligned}$$

Applying Lemma 7 yields:

$$\begin{aligned} &\leq 4d^2 \left( \frac{en}{4m(m-n)t} \right)^t 2e \left( \frac{2t}{c_3 e} \right)^{2t} \\ &\leq 8ed^2 \left( \frac{t}{em(d-1)c_3^2} \right)^t \end{aligned}$$

Since  $\gamma(C_e)$  is stochastically dominated by  $\gamma(H)$ ,

$$\Pr(\gamma(C_e) \geq t) \leq 8ed^2 \left( \frac{t}{e(d-1)c_3^2 m} \right)^t$$

$\square$

This immediately implies the following corollary:

**Corollary 2.** *For any component  $C_e$  found by the Edge Component Search Algorithm,*

$$\Pr(\mathbf{ex}(C_e) \geq s) \leq c_4 \left( \frac{c_5(s+1)}{m} \right)^{s+1}$$

where  $c_4 = 8\epsilon d^2$  and  $c_5 = \frac{1}{\epsilon(d-1)c_3^2} = \frac{d^2(d+1)^2}{\epsilon(d-1)^5}$ .

Note that this bound not only applies to the first component found, but to every component found.

Let  $C'_e$  be the component containing  $e$  if  $e$  is the first edge found in the Edge Component Search Algorithm and let  $C'_e$  be an empty component otherwise.  $\mathbf{ex}(C'_e) = \mathbf{ex}(C_e)$  if  $e$  is the first edge found in  $C_e$ , and  $\mathbf{ex}(C'_e) = 0$  otherwise. In either case the bound of Corollary 2 applies to  $\mathbf{ex}(C'_e)$ .

We will need the following Lemma (proven in section 6).

**Lemma 9.** *Let  $U(s, q)$  be the set of sequences of positive integers, where  $T \in U$  if and only if  $|T| = q$  and  $\sum_{1 \leq i \leq q} T_i = s$ , where  $s \geq q \geq 1$ . Then  $\sum_{T \in U(s, q)} \prod_{1 \leq i \leq q} (T_i + 1)^{T_i+1} \leq 0.89^{q-1} (s+1)^{s+1}$*

We can now bound the excess of the entire graph,

$$\begin{aligned} \Pr(\mathbf{ex}(G) \geq s) &= \Pr\left(\sum_e \mathbf{ex}(C'_e) \geq s\right) \\ &\leq \sum_{\substack{j_1, \dots, j_n \\ \sum_i j_i = s}} \Pr(\wedge_i \mathbf{ex}(C'_{e_i}) \geq j_i) \\ &\leq \sum_{q=1}^s \sum_{\substack{j_1, \dots, j_n \\ \sum_i j_i = s \\ |\{j_i: j_i \geq 1\}| = q}} \prod_{\{j_i: j_i \geq 1\}} c_4 \left( \frac{c_5(j_i + 1)}{m} \right)^{j_i+1} \\ &\leq \sum_{q=1}^s \sum_{\substack{R \subseteq \{1, \dots, n\} \\ |R| = q}} \sum_{\substack{j_1, \dots, j_q \\ \sum_i j_i = s \\ j_i \geq 1}} \prod_{i=1}^q c_4 \left( \frac{c_5(j_i + 1)}{m} \right)^{j_i+1} \\ &\leq \sum_{q=1}^s \binom{n}{q} \sum_{T \in U(s, q)} \prod_{i=1}^q c_4 \left( \frac{c_5(T_i + 1)}{m} \right)^{T_i+1} \\ &\leq \sum_{q=1}^s \binom{n}{q} c_4^q \left( \frac{c_5}{m} \right)^{s+q} \sum_{T \in U(s, q)} \prod_{i=1}^q (T_i + 1)^{T_i+1} \end{aligned}$$

Apply Lemma 9:

$$\leq \sum_{q=1}^s \binom{n}{q} c_4^q \left( \frac{c_5}{m} \right)^{s+q} 0.89^{q-1} (s+1)^{s+1}$$

$$\begin{aligned}
&\leq \frac{1}{0.89}(s+1)^{s+1} \left(\frac{c_5}{m}\right)^s \sum_{q=1}^s \left(\frac{0.89ec_4c_5}{dq}\right)^q \\
&\leq \frac{(s+1)e}{0.89} \left(\frac{c_5s}{m}\right)^s \sum_{q=1}^s \left(\frac{0.89ec_4c_5}{dq}\right)^q
\end{aligned}$$

Applying Lemma 10:

$$\begin{aligned}
&\leq \frac{(s+1)e}{0.89} \left(\frac{c_5s}{m}\right)^s \frac{2e0.89c_4c_5}{d} e^{\frac{0.89c_4c_5}{d}} \\
&\leq \frac{(s+1)16e^2d^3(d+1)^2}{(d-1)^5} e^{\frac{7.2d^3(d+1)^2}{(d-1)^5}} \left(\frac{c_5s}{m}\right)^s
\end{aligned}$$

This completes the proof of Theorem 2.

## 6 Additional Lemmas

This section contains proofs of (more tedious) lemmas.

We begin by showing a useful lemma for inequalities with exponentials.

**Lemma 2.**

$$a_2b_1 \leq a_1b_2 \Rightarrow \left(1 - \frac{1}{a_1+1}\right)^{b_1} \geq \left(1 - \frac{1}{a_2}\right)^{b_2}$$

*Proof.*

$$\begin{aligned}
&\left(1 - \frac{1}{a_1+1}\right)^{b_1} \geq \left(1 - \frac{1}{a_2}\right)^{b_2} \Leftrightarrow \\
&b_1 \ln \left(1 - \frac{1}{a_1+1}\right) \geq b_2 \ln \left(1 - \frac{1}{a_2}\right) \Leftrightarrow \\
&b_1 \left(1 - \frac{1}{1 - \frac{1}{a_1+1}}\right) \geq b_2 \left(-\frac{1}{a_2}\right) \Leftrightarrow \\
&b_1 \left(1 - \frac{a_1+1}{a_1}\right) \geq -\frac{b_2}{a_2} \Leftrightarrow \\
&-\frac{b_1}{a_1} \geq -\frac{b_2}{a_2} \Leftrightarrow \\
&a_2b_1 \leq a_1b_2
\end{aligned}$$

□

**Lemma 3.**  $\text{Bin}(n - q, \frac{1}{m-q+1})$  is stochastically dominated by  $\text{Bin}(n, \frac{1}{m})$  when  $m \geq n$ .

*Proof.* Now  $\text{Bin}(n_1, p_1)$  is stochastically dominated by  $\text{Bin}(n_2, p_2)$  if and only if  $n_1 \leq n_2$  and  $(1 - p_1)^{n_1} \geq (1 - p_2)^{n_2}$  [KM10]. Set  $n_1 = n - q$ ,  $n_2 = n$ ,  $p_1 = \frac{1}{m-q+1}$ ,  $p_2 = \frac{1}{m}$ . Clearly  $n_1 \leq n_2$ . Observe that

$$\left(1 - \frac{1}{m-q+1}\right)^{n-q} \geq \left(1 - \frac{1}{m}\right)^n \Leftarrow$$

Applying Lemma 2

$$\begin{aligned} (n-q)m &\leq n(m-q) \Leftrightarrow \\ -qm &\leq -qn \Leftrightarrow \\ m &\geq n \end{aligned}$$

The last statement is true, so the condition  $(1 - p_1)^{n_1} \geq (1 - p_2)^{n_2}$  holds,  $\square$

**Lemma 4.** *If integers  $m, n, u$  satisfy  $m > n \geq u \geq 1$  then  $\text{Bin}(n - u, \frac{1}{(m-u)^2})$  is stochastically dominated by  $\text{Bin}(n, \frac{1}{m(m-u)})$ .*

*Proof.*  $\text{Bin}(n_1, p_1)$  is stochastically dominated by  $\text{Bin}(n_2, p_2)$  if and only if  $n_1 \leq n_2$  and  $(1 - p_1)^{n_1} \geq (1 - p_2)^{n_2}$  [KM10]. Set  $n_1 = n - u$ ,  $n_2 = n$ ,  $p_1 = \frac{1}{(m-u)^2}$  and  $p_2 = \frac{1}{m(m-u)}$ . Clearly  $n_1 \leq n_2$ . Now

$$\left(1 - \frac{1}{(m-u)^2}\right)^{n-u} \geq \left(1 - \frac{1}{m(m-u)}\right)^n \Leftarrow$$

Applying Lemma 2

$$\begin{aligned} (n-u)m(m-u) &\leq n((m-u)^2 - 1) \\ (n-u)m &\leq n\left(m-u - \frac{1}{m-u}\right) \Leftrightarrow \\ um &\geq n\left(u + \frac{1}{m-u}\right) \Leftrightarrow \\ \frac{m}{n} &\geq \frac{u + \frac{1}{m-u}}{u} \Leftrightarrow \\ \frac{m-n}{n} &\geq \frac{1}{u(m-u)} \Leftrightarrow \\ \frac{n}{m-n} &\leq u(m-u) \end{aligned}$$

The last statement holds since  $u(m-u) \geq m-1 \geq n \geq \frac{n}{m-n}$ . Hence  $(1 - p_1)^{n_1} \geq (1 - p_2)^{n_2}$  as required, so  $\text{Bin}(n - u, \frac{1}{(m-u)^2})$  is stochastically dominated by  $\text{Bin}(n, \frac{1}{m(m-u)})$ .  $\square$

**Theorem 7.** *In the Edge Component Search Algorithm, the number of additional edges between a partially opened vertex pair is stochastically dominated by the number of edges between an unopened vertex pair.*

*Proof.* The case of the pair  $v_1, v_2$  is special because the initial edge  $e$  that was found between these was discovered by selecting a random undiscovered edge, rather than requesting information about the pair  $v_1, v_2$ . Therefore, the occurrence of  $e$  between  $v_1$  and  $v_2$  does not affect the distribution of other edges. Hence the remaining edges between  $v_1$  and  $v_2$  will actually be distributed exactly the same as between any unopened vertex pair.

For the remaining partially opened vertex pairs, we will prove a slightly different statement, which implies the one above. Let  $A$  be the number of additional edges between the partially opened pair and  $B$  be the number between the unopened pair. We show that if  $i < j$ , the probability that  $A = i$  and  $B = j$  is less than the probability that  $A = j$  and  $B = i$ .

Let there be some state,  $S$ , observed on the remainder of the system. By Bayes:

$$\frac{\Pr(B = i \wedge A = j + 1 | S)}{\Pr(B = j \wedge A = i + 1 | S)} = \frac{\Pr(B = i \wedge A = j + 1) \Pr(S | B = i \wedge A = j + 1)}{\Pr(B = j \wedge A = i + 1) \Pr(S | B = j \wedge A = i + 1)}$$

Now  $\Pr(S | B = i \wedge A = j + 1) = \Pr(S | B = j \wedge A = i + 1)$ , since in both cases  $i + j + 1$  edges will have been used between the two vertex-pairs. Therefore, the probability above is simply

$$\frac{\Pr(B = i \wedge A = j + 1)}{\Pr(B = j \wedge A = i + 1)}$$

Note that this statement is true regardless of what the state  $S$  is (as long as it is possible), so we can consider  $S$  to be all information learned about the assignment of edges from the beginning of the Edge Component Search Algorithm.

So, if there are initially  $b$  vertex pairs and  $n$  edges:

$$\begin{aligned} \frac{\Pr(B = i \wedge A = j + 1 | S)}{\Pr(B = j \wedge A = i + 1 | S)} &= \frac{\Pr(B = i \wedge A = j + 1)}{\Pr(B = j \wedge A = i + 1)} \\ &= \frac{\Pr(B = i | A = j + 1) \Pr(A = j + 1)}{\Pr(B = j | A = i + 1) \Pr(A = i + 1)} \\ &= \frac{\binom{n-j-1}{i} \left(\frac{1}{b-1}\right)^i \left(1 - \frac{1}{b-1}\right)^{n-j-1-i} \binom{n}{j+1} \left(\frac{1}{b}\right)^{j+1} \left(1 - \frac{1}{b}\right)^{n-j-1}}{\binom{n-i-1}{j} \left(\frac{1}{b-1}\right)^j \left(1 - \frac{1}{b-1}\right)^{n-i-1-j} \binom{n}{i+1} \left(\frac{1}{b}\right)^{i+1} \left(1 - \frac{1}{b}\right)^{n-i-1}} \\ &= \frac{\binom{n-j-1}{i} \binom{n}{j+1}}{\binom{n-i-1}{j} \binom{n}{i+1}} \end{aligned}$$



$$= \frac{i+1}{j+1} < 1$$

Therefore,  $\Pr(B = i \wedge A = j+1|S) < \Pr(B = j \wedge A = i+1|S)$ , when  $i < j$ , which implies that after observation of the system  $S$ ,  $A - 1$  is stochastically dominated by  $B$ . Therefore, the number of additional edges between a partially opened pair is stochastically dominated by the number of edges between an unopened pair.  $\square$

**Lemma 5.** *Let  $X$  be the sum of independent Bernoulli variables, with mean  $\mu$ . A basic form of the Chernoff Bound for any  $\delta > 0$  is as follows:*

$$\Pr(X \geq (1 + \delta)\mu) \leq \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^\mu$$

This implies the following looser bound:

**Lemma 6.** *For any non-negative integer  $t$ ,*

$$\Pr(X \geq t) \leq \left( \frac{e\mu}{t} \right)^t$$

*Proof.* For  $t \leq \mu$ ,  $\frac{e\mu}{t} \geq e$ , so  $\left(\frac{e\mu}{t}\right)^t \geq 1$ , so the statement holds as the probability cannot be more than 1. For  $t > \mu$  we can view  $t = (1 + \delta)\mu$  for some  $\delta > 0$ .

$$\begin{aligned} \Pr(X \geq (1 + \delta)\mu) &\leq \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^\mu \\ &\leq \left( \frac{e}{(1 + \delta)} \right)^{(1+\delta)\mu} e^{-\mu} \\ &\leq \left( \frac{e\mu}{(1 + \delta)\mu} \right)^{(1+\delta)\mu} e^{-\mu} \\ &\leq \left( \frac{e\mu}{t} \right)^t \end{aligned}$$

$\square$

**Lemma 7.** *Let  $t \geq 1$ ,  $c_3 \in (0, 1)$ . Then:*

$$\sum_{k=1}^{\infty} k^{2t-1} e^{-c_3 k} \leq 2e \left( \frac{2t}{c_3 e} \right)^{2t}$$

*Proof.* It is possible to approximate a summation with an integral, using the same methods as Riemann sums but in reverse. Let  $f(x)$  be a continuous function that is monotonically increasing until a maximum point  $x_{max}$ , after which  $x$

is monotonically decreasing. Let  $x' = \lfloor x_{max} \rfloor$ . Let  $h(x) = \min(f(\lfloor x \rfloor), f(\lfloor x \rfloor + 1))$ . Let us observe how  $\sum_{x=a}^b h(x)$  approximates  $\int_a^{b+1} f(x)dx$ .

Observe that for any integer  $a$ ,  $h(x)$  is the same for all  $x \in [a, a+1)$ . Since  $f(x)$  has no local minima and is continuous, the minimum value of  $f(x)$  over the range  $[a, a+1)$  is either at  $f(a)$  or  $f(a+1)$ . Therefore  $f(x) \geq \min(f(a), f(a+1)) = h(x)$  for  $x \in [a, a+1)$ . Since this applies to the interval  $[a, a+1)$  for any integer  $a$ ,  $h(x) \leq f(x)$  for all  $x$ . Hence, for any integers  $a$  and  $b$ ,  $\sum_a^b h(x) = \int_a^{b+1} h(x) \leq \int_a^{b+1} f(x)dx$ .

$$\begin{aligned} \int_a^{b+1} f(x)dx &\geq \sum_a^b h(x) \\ &\geq \sum_a^{x'-1} f(x) + \min(f(x'), f(x'+1)) + \sum_{x'+1}^b f(x+1) \\ &\geq \left( \sum_a^{b+1} f(x) \right) - \max(f(x'), f(x'+1)) \\ &\geq \left( \sum_a^{b+1} f(x) \right) - f(x_{max}) \end{aligned}$$

Hence:

$$\sum_a^b f(x) \leq \int_a^b f(x)dx + f(x_{max})$$

Let  $f(x) = x^{2t-1}e^{-c_3x}$  where  $t \geq 1$  and  $0 < c_3 < 1$ . First we need to show that it is a function that is monotonically increasing, then monotonically decreasing.

$$f'(x) = (2t-1)x^{2t-2}e^{-c_3x} - c_3x^{2t-1}e^{-c_3x} = x^{2t-2}e^{-c_3x}(2t-1-c_3x)$$

Observe that  $x^{2t-2}$  and  $e^{-c_3x}$  are both positive. Therefore  $f'(x)$  will be positive when  $x < \frac{2t-1}{c_3}$ ,  $f'(x) = 0$  at  $x = \frac{2t-1}{c_3}$  and will be negative when  $x > \frac{2t-1}{c_3}$ . Therefore this function is monotonically increasing, then monotonically decreasing, as required, with  $x_{max} = \frac{2t-1}{c_3}$ . We can easily calculate:

$$\begin{aligned} f(x_{max}) &= \left( \frac{2t-1}{c_3} \right)^{2t-1} e^{-(2t-1)} \\ &= \left( \frac{2t-1}{c_3 e} \right)^{2t-1} \end{aligned}$$

Hence the inequality applies to the sum and

$$\sum_{k=1}^{\infty} k^{2t-1} e^{-c_3k} \leq \int_1^{\infty} x^{2t-1} e^{-c_3x} dx + \left( \frac{2t-1}{c_3 e} \right)^{2t-1}$$

$$\leq \int_0^\infty x^{2t-1} e^{-c_3 x} dx + \left( \frac{2t-1}{c_3 e} \right)^{2t-1}$$

By a standard integral identity,  $\int_0^\infty x^{2t-1} e^{-c_3 x} dx = \frac{(2t-1)!}{c_3^{2t}}$ . Furthermore, a factorial approximation shows that  $(2t-1)! \leq (2t)^{2t} e^{-(2t-1)}$ . Hence

$$\begin{aligned} \sum_{k=1}^\infty k^{2t-1} e^{-c_3 k} &\leq \left( \frac{2t}{c_3} \right)^{2t} e^{-(2t-1)} + \left( \frac{2t-1}{c_3 e} \right)^{2t-1} \\ &\leq \left( \frac{2t}{c_3} \right)^{2t} e^{-(2t-1)} + \left( \frac{2t}{c_3} \right)^{2t-1} e^{-(2t-1)} \end{aligned}$$

Recalling that  $0 < c_3 < 1$ , so  $\frac{2t}{c_3} > 1$

$$\begin{aligned} &\leq 2 \left( \frac{2t}{c_3} \right)^{2t} e^{-(2t-1)} \\ &\leq 2e \left( \frac{2t}{c_3 e} \right)^{2t} \end{aligned}$$

□

**Lemma 8.** For all integers  $s \geq 2$ ,  $\sum_{a=1}^{s-1} (a+1)^{a+1} (s+1-a)^{s+1-a} \leq 0.89(s+1)^{s+1}$

*Proof.* By calculation, this is true for  $s \in \{2, 3, 4, 5, 6, 7, 8\}$ , for which the left-hand side values are, respectively  $\{16, 216, 2777, 38824, 607534, 10707768, 212342547\}$  and the right-hand side values are respectively  $\{24.03, 227.84, 2781.25, 41523.84, 732953.27, 14931722.24, 344804235.2\}$ . For  $s > 8$  we prove by induction.

Given that it holds true for  $s \leq 8$ , let us show it holds true for  $s+1$ .

$$\begin{aligned} &\sum_{a=1}^s (a+1)^{a+1} (s+2-a)^{s+2-a} \\ &\leq \sum_{a=1}^t (a+1)^{a+1} (s+2-a)^{s+2-a} + \sum_{a=t+1}^{s-1} (a+1)^{a+1} (s+2-a)^{s+2-a} + (s+1)^{s+1} 2^2 \\ &\leq \sum_{a=1}^t (a+1)^{a+1} (s+2-a)^{s+2-a} + (s+1-t)e \sum_{a=t+1}^{s-1} (a+1)^{a+1} (s+1-a)^{s+1-a} + (s+1)^{s+1} 2^2 \\ &\leq \sum_{a=1}^t (a+1)^{a+1} (s+2-a)^{s+2-a} + (s+1-t)e \end{aligned}$$

$$\left( \sum_{a=1}^{s-1} (a+1)^{a+1} (s+1-a)^{s+1-a} - \sum_{a=1}^t (a+1)^{a+1} (s-a+1)^{s-a+1} \right) + (s+1)^{s+1} 2^2$$

Applying the inductive hypothesis yields:

$$\begin{aligned} &\leq e \sum_{a=1}^t (a+1)^{a+1} (s+2-a)(s+1-a)^{s+1-a} + (s+1-t)e0.89(s+1)^{s+1} \\ &\quad - e(s+1-t) \sum_{a=1}^t (a+1)^{a+1} (s-a+1)^{s-a+1} + (s+1)^{s+1} 2^2 \\ &\leq e \sum_{a=1}^t (a+1)^{a+1} (s-a+1)^{s-a+1} ((s+2-a) - (s+1-t)) + 0.89(s+2)^{s+2} \\ &\quad - te0.89(s+1)^{s+1} + (s+1)^{s+1} 2^2 \\ &\leq 0.89(s+2)^{s+2} + 2^2(s+1)^{s+1} + e \sum_{a=1}^t (a+1)^{a+1} (s-a+1)^{s-a+1} (t+1-a) - te0.89(s+1)^{s+1} \end{aligned}$$

Setting  $t = 3$  yields:

$$\begin{aligned} &\leq 0.89(s+2)^{s+2} + 2^2(s+1)^{s+1} + e2^2s^s3 + e3^3(s-1)^{s-1}2 + e4^4(s-2)^{s-2} - 3e0.89(s+1)^{s+1} \\ &\leq 0.89(s+2)^{s+2} + (s+1)^{s+1} \left( 2^2 + \frac{12}{s} + \frac{54}{es(s-1)} + \frac{256}{e^2s(s-1)(s-2)} - 3e0.89 \right) \end{aligned}$$

For  $s \geq 8$ , the term  $2^2 + \frac{12}{s} + \frac{54}{es(s-1)} + \frac{256}{e^2s(s-1)(s-2)} \leq 5.5$ . Since  $3e0.89 > 5.5$  the inequality simplifies to:

$$\sum_{a=1}^s (a+1)^{a+1} (s+2-a)^{s+2-a} \leq 0.89(s+2)^{s+2}$$

Since it holds true up to  $s = 2, \dots, 8$  by inspection, and holds true for  $s \geq 8$  by induction, the statement is true for all  $s \geq 2$ .  $\square$

**Lemma 9.** Let  $U(s, q)$  be the set of sequences of positive integers, where  $T \in U$  if and only if  $|T| = q$  and  $\sum_{1 \leq i \leq q} T_i = s$ , where  $s \geq q \geq 1$ . Then  $\sum_{T \in U(s, q)} \prod_{1 \leq i \leq q} (T_i + 1)^{T_i+1} \leq 0.89^{q-1} (s+1)^{s+1}$

*Proof.* We proceed by induction on the length of the sequences. For  $q = 1$ ,  $U$  contains a single sequence  $T$  with  $T_1 = s$ . Then  $\sum_{T \in U(s, q)} \prod_{1 \leq i \leq q} (T_i + 1)^{T_i+1} = (s+1)^{s+1} = 0.89^0 (s+1)^{s+1}$ .

Assume that the theorem holds for all sequences of length  $q \geq 1$ . We will show that it also holds for all sequences of length  $q + 1$ .

$$\sum_{T \in U(s, (q+1))} \prod_{1 \leq i \leq q+1} (T_i + 1)^{T_i+1} \leq \sum_{T_1=1}^{s-q} (T_1 + 1)^{T_1+1} \sum_{T' \in U((s-T_1), q)} \prod_{1 \leq i \leq q} (T'_i + 1)^{T'_i+1}$$

Applying our inductive hypothesis gives:

$$\begin{aligned} &\leq \sum_{a=1}^{s-1} (a+1)^{a+1} 0.89^{q-1} (s-a+1)^{s-a+1} \\ &\leq 0.89^{q-1} \sum_{a=1}^{s-1} (a+1)^{a+1} (s-a+1)^{s-a+1} \end{aligned}$$

Using Lemma 8

$$\leq 0.89^q (s+2)^{s+2}$$

□

**Lemma 10.** For  $a > 0$ ,

$$\sum_{x=1}^{\infty} \left(\frac{a}{x}\right)^x \leq 2ae^{\frac{a}{e}}$$

*Proof.* Over the positive reals the function  $f(x) = \left(\frac{a}{x}\right)^x$  is maximized at  $x = \frac{a}{e}$ , for which it has value  $e^{\frac{a}{e}}$ . Therefore:

$$\begin{aligned} \sum_{x=1}^{\infty} \left(\frac{a}{x}\right)^x &\leq \sum_{x=1}^{2a-1} \left(\frac{a}{x}\right)^x + \sum_{x=2a}^{\infty} \left(\frac{a}{x}\right)^x \\ &\leq \sum_{x=1}^{2a-1} \left(e^{\frac{a}{e}}\right) + \sum_{x=2a}^{\infty} \left(\frac{1}{2}\right)^x \\ &\leq e^{\frac{a}{e}}(2a-1) + 1 \\ &\leq 2ae^{\frac{a}{e}} \end{aligned}$$

□

## 7 A Lower Bound

We now show the following lower bound on the number of elements  $n$  in terms of the security parameter  $N$ , such that cuckoo hashing with a stash can fail with negligible probability in  $N$ . For consistency with other parts of the paper, we use the 2-table construction but this can easily be adapted to other constructions.

**Theorem 13.** *If  $n = \mathcal{O}(\log(N))$  and  $n - s = \Omega(n)$  then it is impossible for a 2-table Cuckoo Hash table to have a negligible build failure probability in  $N$ .*

*Proof.* Since  $n - s = \Omega(n)$ , it follows that  $n - s \geq c_0 n$  for sufficiently large  $n$  where the constant  $c_0$  satisfies  $0 < c_0 \leq 1$ . Therefore:

$$\begin{aligned} \frac{n - s}{n} &\geq c_0 \\ \frac{n - s - 2}{n} &\geq c_0 - \frac{2}{n} \\ \frac{n - s - 2}{n} &\geq \frac{c_0}{2} \text{ when } n \geq \frac{4}{c_0} \\ \frac{n - s - 2}{n} &\geq c_1 \text{ for constant } c_1 \text{ satisfying } 0 < c_1 \leq \frac{1}{2} \end{aligned}$$

Since  $n = \mathcal{O}(\log(N))$ , there is some constant  $c_2$  such that  $n \leq c_2 \log(N)$  (for sufficiently large  $n$ ).

Let  $m = dn$  be the size of each table.

If all  $n$  items are hashed to the first  $\lceil \frac{n-s-2}{2} \rceil$  locations in both tables, then  $2\lceil \frac{n-s-2}{2} \rceil \leq n - s - 1$  items can be stored in the table, and  $s$  items can be stored in the stash, but 1 item will not be able to be stored at all, so the build fails.

The probability that all  $n$  items are stored in the first  $\lceil \frac{n-s-2}{2} \rceil$  locations in both tables is at least:

$$\begin{aligned} \left(\frac{n - s - 2}{2dn}\right)^{2n} &\geq \left(\frac{c_1}{2d}\right)^{2c_2 \log(N)} \\ &\geq N^{2c_2 \log(\frac{c_1}{2d})} \end{aligned}$$

This is non-negligible in  $N$ . Therefore the probability of a build failure is non-negligible. □

This immediately implies the contrapositive:

**Corollary 3.** *Cuckoo Hashing with a stash requires  $n - s = o(n)$  or  $n = \omega(\log(N))$  in order to succeed with failure negligible in  $N$ .*

The case that  $n - s = o(n)$  is very unnatural—it implies that a sub-constant number of elements are stored in the table, at which point the Cuckoo table is not providing much use. Thus, in any realistic setting where Cuckoo tables are used, it is necessary that  $n = \omega(\log(N))$ . This provides the lower bound for  $n$  in terms of  $N$  such that Cuckoo Hashing with a stash has a negligible probability of failure.

## 8 Acknowledgments

This research was sponsored in part by ONR grant (N00014-15-1-2750) “Syn-Crypt: Automated Synthesis of Cryptographic Constructions”.

## References

- [ADW14] Martin Aumüller, Martin Dietzfelbinger, and Philipp Woelfel. Explicit and efficient hash families suffice for cuckoo hashing with a stash. *Algorithmica*, 70(3):428–456, 2014.
- [Aum10] Martin Aumüller. *An alternative analysis of cuckoo hashing with a stash and realistic hash functions*. PhD thesis, Diplomarbeit, Technische Universität Ilmenau, 2010.
- [DK12] Michael Drmota and Reinhard Kutzelnigg. A precise analysis of cuckoo hashing. *ACM Transactions on Algorithms (TALG)*, 8(2):1–36, 2012.
- [Dwa69] Meyer Dwa. The total progeny in a branching process and a related random walk. *Journal of Applied Probability*, 6(3):682–686, 1969.
- [GM11] Michael T Goodrich and Michael Mitzenmacher. Privacy-preserving access of outsourced data via oblivious RAM simulation. In *ICALP*, pages 576–587. Springer, 2011.
- [GMOT12] Michael T Goodrich, Michael Mitzenmacher, Olga Ohrimenko, and Roberto Tamassia. Privacy-preserving group data access via stateless oblivious RAM simulation. In *SODA*, pages 157–167. SIAM, 2012.
- [KLO12] Eyal Kushilevitz, Steve Lu, and Rafail Ostrovsky. On the (in) security of hash-based oblivious RAM and a new balancing scheme. In *SODA*, pages 143–156. SIAM, 2012.
- [KM10] Achim Klenke and Lutz Mattner. Stochastic ordering of classical discrete distributions. *Advances in Applied probability*, 42(2):392–410, 2010.
- [KMW09] Adam Kirsch, Michael Mitzenmacher, and Udi Wieder. More robust hashing: Cuckoo hashing with a stash. *SIAM Journal on Computing*, 39(4):1543–1561, 2009.
- [Kut10] Reinhard Kutzelnigg. A further analysis of cuckoo hashing with a stash and random graphs of excess  $r$ . *Discrete Mathematics and Theoretical Computer Science*, 12(3):81–101, 2010.
- [Pit98] Jim Pitman. Enumerations of trees and forests related to branching processes and random walks. *Microsurveys in discrete probability*, 41:163–180, 1998.
- [PR01] Rasmus Pagh and Flemming Friche Rodler. Cuckoo hashing. In *ESA*, pages 121–133. Springer, 2001.



- [PR10] Benny Pinkas and Tzachy Reinman. Oblivious RAM revisited. In *CRYPTO*, pages 502–519. Springer, 2010.
- [PSSZ15] Benny Pinkas, Thomas Schneider, Gil Segev, and Michael Zohner. Phasing: Private set intersection using permutation-based hashing. In *24th USENIX Security Symposium (USENIX Security 15)*, pages 515–530, 2015.
- [PSWW18a] Benny Pinkas, Thomas Schneider, Christian Weinert, and Udi Wieder. Efficient circuit-based psi via cuckoo hashing. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 125–157. Springer, 2018.
- [PSWW18b] Benny Pinkas, Thomas Schneider, Christian Weinert, and Udi Wieder. Efficient circuit-based psi via cuckoo hashing. *Cryptology ePrint Archive*, 2018.
- [Wie16] Udi Wieder. Hashing, load balancing and multiple choice draft, 2016.