

A More Elaborate Computation Model for Analyzing Balloon Hashing

Cui Hongrui, Sha Jinrui

December 11, 2017

1 An improved big-spread lemma

Room for improvement The *big spread lemma* plays an important role in the security proof of balloon hashing algorithm [1], however, the theorem require that the subset being considered, V' have the same size of the pebbles available by the adversary. This is also reflected in Lemma 30, which bound the probability of the graph being too “dense”. Recall in Lemma 30, for $\delta = 3$, every sandwich graph is an $(m, m, n/16)$ -consecutively-avoiding graph, for $n_0 > 16$ and all $n_0 \leq m < n/6$ except with probability $P_{consec}(n, d, n_0) \leq 2 * 8 \cdot d \cdot n \cdot 2^{-n_0/2}$; And for $\delta = 7$, such stack of sandwich graph is $(n/64, n/32)$ -everywhere-avoiding graph when $n > 2^{11}$, except with probability $P_{every} \leq 128 \cdot d \cdot 2^{-n/50}$. Note that in the later part of the lemma, the size of V' is defined as $m = \frac{n}{2^{\omega+1}}$, and thus m grows with n .

This limitation that the subset of V begin fixed by the number of pebbles of the adversary is somehow inconvenient, in that to make the probability of bad event small, one need the number of pebbles available grow, which is the contrary of natural instinct, as the more pebble the adversary acquires, the easier the attack (re-computation) should be. Although the result in [1] is absolutely correct, as we can view the growth of m as the result of the growth of n , an if the subset V' being considered too small, the variables sampled independently uniform at random from the last layer would have a bigger probability to being not so well-spread. Still, we want to separate the size of subset (m) and the number of pebbles ($|V'|$).

For the sake of clarity, and also because of the nature of an internal manuscript, I will replicate the content and the proof of big spread lemma.

Lemma 1 (Big Spread Lemma Variation). *For all positive integer $\delta \geq 3$, $\omega \geq 2$, n_0 , and n , and for all positive integers m such that $2^\omega < m < 2^{2-\omega+\frac{2\omega}{\delta e}}$, a list of δm elements sampled independently and uniformly at random from $1, \dots, n$ is an $(n_0, n/2^\omega)$ -well-spread set with probability at least $1 - 2^{(1-\omega)\frac{\delta}{2}m} \cdot (2^{\omega m} + 2^{\omega+1} 2^{\omega n_0})$*

Proof. The Strategy to bound the bad event that such property does not exist is the same as that used in [1]. Let $R = (R_1, \dots, R_{\delta m})$ be integers sampled independently and uniformly at random from $1, \dots, n$, we want to prove that for all subset $S \subseteq R$ of size at most n_0 , $spread_S(R) \geq n/2^\omega$. To do so, we first define a bad event B, and then show that bounding $\Pr[B]$ is enough to prove the lemma, then bound $\Pr[B]$ to complete the lemma.

The Bad Event B Write the integers in R in non-decreasing order as $X_1, \dots, X_{\delta m}$, then define $X_0 = 0, X_{\delta m+1} = n$. Let bad event B be the event that there exists a set $S' \subseteq X_1, \dots, X_{\delta m+1}$ of size at most $(n_0 + 1)$, such that $\sum_{X_i \in S'} X_i - X_{i-1} \geq (1 - 2^{-\omega})n$.

Whenever there exists a set $S \subseteq R$ of size at most n_0 that cause $\text{spread}_S(R) < n/2^\omega$, then bad event B must occur. Assuming that such a bad set S exists, construct a set $S' = S \cup X_{\delta m+1}$ of size at most $n_0 + 1$, Then we compute

$$\sum_{X_i \in S'} (X_i - X_{i-1}) = n - X_{\delta m} + \sum_{X_i \in S} (X_i - X_{i-1}) = n - \sum_{X_i \notin S} (X_i - X_{i-1}) \quad (1)$$

The last inequality holds because $X_{\delta m} = \sum_{X_i \in R} (X_i - X_{i-1})$. Now that $\text{Spread}_S(R) < n/2^\omega$, so we have bad event B occurs. Thus, $\Pr[\text{Spread}_S(R)] \leq \Pr[B]$, and therefore bounding a upper bound of bad event B is sufficient to prove the lemma.

Strategy to bound $\Pr[B]$ Let D be a random variable denoting the number of distinct integers in the list of random integers R . For any fixed integer $d^* \in \{1, \dots, n\}$, we can write:

$$\Pr[B] = \Pr[B|D < d^*] \cdot \Pr[D < d^*] + \Pr[B|D \geq d^*] \cdot \Pr[D \geq d^*] \quad (2)$$

$$\leq \Pr[D < d^*] + \Pr[B, D \geq d^*] \quad (3)$$

In the following proof, we take $d^* = \delta m/2$.

Bounding $\Pr[D < d^*]$ The probability that this event occurs is at most the probability when we throw δm balls into n bins, all the balls fall into a set of $\delta m/2$ bins.

$$\Pr[D < d^*] \leq \binom{n}{d^*} \left(\frac{d^*}{n}\right)^{\delta m} \leq \left(\frac{n \cdot e}{d^*}\right)^{d^*} \left(\frac{d^*}{n}\right)^{\delta m} = \left(\frac{d^*}{n}\right)^{\delta m - d^*} e^{d^*} \quad (4)$$

By the hypothesis of the lemma, $m < 2^{(1-\omega+\frac{2\omega}{\delta})\frac{2n}{\delta e}}$. Therefore, we have

$$\Pr[D < d^*] \leq \left(\frac{\delta m e}{2n}\right)^{\delta m/2} \leq 2^{[(1-\omega)\frac{\delta}{2} + \omega]m} \quad (5)$$

Bounding $\Pr[B, D \geq d^*]$ First I would like to re-state the lemma 8 in the original work [1], because this lemma is useful in the proof of big spread lemma.

Lemma 2. *Let (R_1, \dots, R_d) be random variables respecting integers sampled uniformly, but without replacement, from $\{1, \dots, n\}$. Write the R s in ascending order as (Y_1, \dots, Y_d) , and let $Y_0 = 0$, and $Y_{d+1} = n$. Next define $L = (L_1, \dots, L_{d+1})$, where $L_i = Y_i - Y_{i-1}$. Then for all functions $f : \mathbb{Z}^{d+1} \rightarrow \{0, 1\}$, and for all permutations π on $d+1$ elements,*

$$\Pr[f(L_1, \dots, L_{d+1}) = 1] = \Pr[f(L_{\pi(1)}, \dots, L_{\pi(d+1)}) = 1]. \quad (6)$$

Particularly, this lemma shows that if we define a function that maps a list of segments (separated from $\{1, \dots, n\}$ by the δm random integers) in to $\{0, 1\}$, such that if bad event B happens, then it is one. Using this lemma, we can define a permutation on those integers that if it appears in a bad set, then we

shift it with the first $n_0 + 1$ integers in the list, and thus we can only focus on the first few elements on the list without changing the result. As lemma 2 require the list to be mutually distinct, we have to transform the random list in to a distinct one first.

Step 1: Bad event B occurs whenever there is a subset $S' \subseteq X = (X_1, \dots, X_{\delta m+1})$ of size at most $(n_0 + 1)$, such that $\sum_{X_i \in S'} (X_i - X_{i-1}) \geq (1 - 2^{-\omega})n$. Whenever bad event B occurs, there is also a subset S'_{dist} of distinct integers that also satisfy the sum equation, and also has the size $(n_0 + 1)$. (This is because if duplicate exists, we can always delete duplicates, without changing the sum, and then add new integers to make the sum even bigger.) thus we can use lemma 2 to simplify calculation.

Step 2: we consider the case there is explicitly d distinct integers in the list R , and then use union bound to deduce the final probability. Write the d distinct integers in R in ascending order as $Y = (Y_1, \dots, Y_d)$. We want to bound the probability that some subset Y'_{dist} of size n_0 is bad. (Note that the only change from the original lemma is that I changed m into n_0 .) Then define L_i as in lemma 2. Let the set of indices $I \subseteq \{1, \dots, d+1\}$ of size $n_0 + 1$ be the set of index of bad subset. That is, $S'_{dist} = \{Y_i | i \in I\}$. Then define the permutation $\pi : \mathbb{Z}^{d+1} \rightarrow \mathbb{Z}^{d+1}$, such that if $i \in I$, then L_i appears in the first $n_0 + 1$ elements of $(L_{\pi(1)}, \dots, L_{\pi(d+1)})$. Define function $f : \mathbb{Z}^{d+1} \rightarrow \{0, 1\}$ such that if its first $n_0 + 1$ arguments sums to at least $(1 - 2^{-\omega})n$, returns 1 and it returns 0 otherwise. Therefore, we have:

$$Pr[\sum_{i \in I} L_i \geq (1 - 2^{-\omega})n | D = d] = Pr[f(L_{\pi(1)}, \dots, L_{\pi(d+1)}) = 1 | D = d] \quad (7)$$

$$= Pr[f(L_1, \dots, L_{d+1}) = 1 | D = d] \quad (8)$$

$$= Pr[(L_1 + \dots + L_{n_0+1}) \geq (1 - 2^{-\omega})n | D = d] \quad (9)$$

Step 3: the event defined above can only happen when the other $d - (n_0 + 1)$ integers falls into the right most $2^{-\omega}n$ bins, and then we have:

$$Pr[(L_1 + \dots + L_{n_0+1}) \geq (1 - 2^{-\omega})n | D = d] \leq \left(\frac{1}{2}\right)^{\omega(d - n_0 - 1)} \quad (10)$$

Above is the probability that a single set is bad, we then sum it over all $\binom{d+1}{n_0+1}$ possible size $n_0 + 1$ subsets, to get:

$$Pr[B | D = d] \leq \binom{d+1}{n_0+1} \left(\frac{1}{2}\right)^{\omega(d - n_0 - 1)} \quad (11)$$

$$\leq 2^{(1-\omega)d + \omega n_0 + \omega + 1} \quad (12)$$

Equipped with $Pr[B | D = d]$, we can then proceed to calculate $Pr[B, D \geq d^*]$. Note when $\omega > 1$, $Pr[B | D = d]$ is non-increasing in d . And therefore we have:

$$Pr[B, D \geq d^*] \leq 2^{(1-\omega)d^* + \omega n_0 + \omega + 1} \cdot \sum_{d=d^*}^{\delta m+1} Pr[D = d] \quad (13)$$

$$\leq 2^{1+\omega} \cdot 2^{(1-\omega)\frac{\delta m}{2} + \omega n_0} \quad (14)$$

Completing the proof By the calculation above, we have

$$Pr[B] \leq Pr[D < d^*] + Pr[B, D \geq d^*] \quad (15)$$

$$\leq 2^{(1-\omega)\frac{\delta m}{2}} \cdot (2^{\omega m} + 2^{1+(n_0+1)\omega}) \quad (16)$$

□

2 Everywhere avoiding property

An analysis of original strategy In the original work [1], the author proves a *consecutive avoiding property* over all possible size of m , and then for m of fixed size ($n/64$), proves a tighter *everywhere avoiding property*. The reason of combining these two property is that the first one allows the number pebbles to be arbitrary within a given range. In particular, one can consider m to have the same number of the pebbles the adversary owns. In order to prove a tighter bound on $S \cdot T$ (like in part b of Lemma 31), one needs to first apply consecutive avoiding property and then use the everywhere avoiding property. What we want to do is to acquire the probability that given at most n_0 pebbles, any m pebbles on a standard sandwich graph is an $(n_0, m, 2/2^\omega)$ avoiding set.

Lemma 3 (Everywhere Avoiding). *Let $G = (U \cup V, E)$, be a δ -random sandwich graph on $2n$ vertices. Then for all positive integer m that satisfy the assumption of lemma 1, and n_0 such that n_0 is no bigger than the upper bound of m , we have that given n_0 pebbles at most, every subset $V' \subseteq V$ of size m , is a $(n_0, n/2^\omega)$ avoiding set with probability at least $1 - (\frac{n}{m})^m \cdot 2^{(1-\omega)\frac{\delta m}{2}} \cdot (2^{\omega m} + 2^{1+(n_0+1)\omega})$*

Proof. Fix the size of n_0 and m , all we need to do is to apply union bound to all possible choice of V' , which has $\binom{n}{m}$ in total. Apply this inequality to simply the bound:

$$\binom{n}{m} = \frac{n \times (n-1) \times \cdots \times (n-m+1)}{m \times (m-1) \times \cdots \times 1} \leq \frac{n^m}{m^m} = \frac{n^m}{m^m} \quad (17)$$

And thus by lemma 1, we have the probability that a single bad event occurs is at most $2^{(1-\omega)\frac{\delta}{2}m} \cdot (2^{\omega m} + 2^{\omega+1}2^{\omega n_0})$, and thus the probability that for all choice of V' , the bad event occurs is at most $(\frac{n}{m})^m \cdot 2^{(1-\omega)\frac{\delta m}{2}} \cdot (2^{\omega m} + 2^{1+(n_0+1)\omega})$. And thus the lemma is proved. □

3 A bound on the pebbling moves of random sandwich graph

Bounding pebbling moves And now for the sake of our own objective, we would like to calculate the bound of re-computation of a subset of vertices on any layer. To do this, we need to first compute the moves needed to pebble a subset of vertices. First define a variable to denote the pebbling moves needed to pebble a subset of m vertices on the condition that except on those m vertices, there can be at most n_0 pebbles on any vertices of the graph.

Definition 1. Let G be a stack of random sandwich graphs. Then define T_r be the pebbling moves needed to pebble a subset V' of vertices on layer r , which has not received any pebble on the beginning of the moves, and at the end of the moves, the topologically last one receives a pebble. The condition is that there can only be at most n_0 pebbles on the graph.

Now that we have defined the amount to calculate, we can use induction reasoning to calculate T_r .

Lemma 4. Given the average case that $T_0 = \frac{m+n}{2}$, we have:

$$T_r \leq -\frac{m}{\frac{n}{m2^\omega} - 1} + \left(\frac{n}{m2^\omega}\right)^r \left(\frac{m+n}{2} + \frac{m}{\frac{n}{m2^\omega} - 1}\right) \quad (18)$$

where $0 \leq r \leq d$, d is the layer of the stack of random sandwich graph.

Proof. By the definition of T_r , we have that $T_r \leq m + \frac{n}{m2^\omega} T_{r-1}$. And thus we have:

$$T_r \leq m + \frac{n}{m2^\omega} T_{r-1} \quad (19)$$

$$T_r + \frac{m}{\frac{n}{m2^\omega} - 1} \leq \frac{n}{m2^\omega} \left(T_{r-1} + \frac{m}{\frac{n}{m2^\omega} - 1}\right) \quad (20)$$

$$(21)$$

If we define $Y_r = T_r + \frac{m}{\frac{n}{m2^\omega} - 1}$, and change the inequality into equality, we can have

$$Y_r = \left(\frac{n}{m2^\omega}\right)^r \cdot Y_0 \quad (22)$$

And thus $T_r \leq -\frac{m}{\frac{n}{m2^\omega} - 1} + \left(\frac{n}{m2^\omega}\right)^r \left(\frac{m+n}{2} + \frac{m}{\frac{n}{m2^\omega} - 1}\right)$. \square

Then we define the lower bound of re-computation needed in the process discussed above.

Definition 2. Let $Recompute(r)$ be the pebbling moves needed in the process to pebble m vertices on layer r . More specifically, the pebbling moves needed to pebble the direct predecessors of those m vertices.

By the definition of $Recompute(r)$ and lemma 4, one can see that $Recompute(r) \geq -m + T_r$, and it grows exponentially with r .

4 A comparison on re-computation and DRAM fetching

Different strategy The pebbling diagram intends to model the computation on ASIC circuits, that is, the computation of underlying cryptographic hashing function, like SHA-256, is relatively very fast¹, and those logic takes the major size of the circuit. The memory section, on the other hand, is relatively small. And thus in the model, we can model the adversary to have small amount of pebbles, and try to bound the pebbling moves of the adversary.

¹Some pipelining circuits can reach the frequency as fast as 1THz

Layer	Ratio
1	0.01
2	0.04
3	0.15
4	0.62
5	2.46

Table 1: The ratio of re-computation over DRAM fetching

However, this model cannot extend to honest users, whose considerable amount of time is used in fetching previously computed blocks from DRAM, and this operation takes constant time. Intuitively, if the adversary, or the manufacturer of ASIC can combine those two advantages, could a more devastating attack be constructed? Especially when DRAM circuits are getting faster and relatively cheaper. On the other hand, when such case occurs, then total computational time depends not only on the speed of hashing, but also on the DRAM latency. In this case, will the increase in hashing power (non-memory hard hashing power) yields a less attractive increase in the throughput, and therefore render the solution of ASIC as a cracking machine useless?

Our preliminary work shows that given contemporary parameter settings, when the layer of the graph gets greater than 4, the DRAM strategy will definitely start to have an advantage.

5 Comparison on common parameter settings

Parameter setting In the following comparison section, we take the block size, which the output length of the underlying cryptographic hashing function to be 256 bit, and the block in one layer of the graph to be $2^{22}2^2$. We take $\omega = 3$, and the ratio of the latency of fetching a DRAM block over the ratio of the time needed to output a block from Random Oracle, $c = 1000$.

Comparison As shown in table 1, given our parameter setting, the ratio of re-computation becomes bigger than 1 when the graph get deeper than 4 layers. However, this is a loose bound. Firstly, the DRAM latency is computed as δmc , where $c = 1000$, this does not take into account the case when the δm predecessors ‘collide’ with each other, that is, when one fetch benefit another pebbling process. Secondly, the re-computation was take the upper bound. Thirdly, an optimal strategy may combine these two means, using fetching for some vertices, and re-computation for others. Nevertheless, this rather loose bound shows that if the layer of the graph gets deeper, pure re-computation will never be a useful solution, no matter what.

Definition 2(Cumulative Pebbling Complexity). Let G be a DAG, $P = (P_0, \dots, P_t)$ be an arbitrary pebbling of G and \mathcal{P} be the set of all complete pebbblings of G . Then the (cumulative) cost of P and then cumulative complexity (CC)

²The authors of Balloon Hashing recommended the space parameter to be 256MB for authentication servers.

of G are defined respectively to be:

$$p - cost(P) := \sum_{i=0}^t |P_i| \quad (23)$$

$$cc(G) := \min\{p - cost(P) : P \in \pi\} \quad (24)$$

Moreover the space complexity (SC) and space/time complexity (STC) are defined as follows:

$$s - cost(P) := \max\{|P_i| : i \in \{0, \dots, t\}\} \quad (25)$$

$$sc(G) := \min\{s - cost(P) : P \in \pi\} \quad (26)$$

$$st - cost(P) := k * \max\{|P_i| : i \in \{0, \dots, t\}\} \quad (27)$$

$$sc(G) := \min\{st - cost(P) : P \in \pi\} \quad (28)$$

Amortized Graph Complexity. Unlike the standard pebbling complexity notions of SC and STC , for the parallel pebbling game, CC is also an amortized complexity notion. In fact it scales additively in the CC of the individual disconnected components of the graph. We prove this property of CC formally in the following lemma. In our model, we need make some changes in the definition CC .

Definition 3(New Cumulative Pebbling Complexity). we suppose that the time which we fetch the block is c times higher than do once hash. At the same time, to make our next analysis become easier we also suppose that we don't make other work when we fetch the block in the DRAM, so based on our this suppose, we define the status when we fetch the block are $Q_i (i = 1, \dots, s)$. So in our model, we have :

$$p' - cost(P) := \sum_{i=0}^t |P_i| + (c - 1) \sum_{i=0}^s |Q_i| \quad (29)$$

$$cc(G) := \min\{p' - cost(P) : P \in \pi\} \quad (30)$$

Definition 4(Heavy Group). For $a \in \mathbb{N}_{< i-\alpha}$, we call $\Gamma_a := \{(\phi, v) \in p : \lfloor \frac{v}{2^\alpha} \rfloor = a\}$ is a group, and a heavy group is one which has $|\Gamma_a| \geq (1 - 3\omega)2^\alpha$

According the article[?], we can get the lemma of the number of the heavy groups:

Lemma 5. There are at least $2^{i-\alpha}/3$ heavy groups

The proof of this lemma is in the article.[?] After give the definition of heavy group, we need to express the complexity of $cc(G_{\omega, \phi, i})$, so we give the Lemma 6:

Lemma 6. Fix any $\phi \in \mathbb{N}, 0 < \omega \leq \frac{4^{-\phi}}{6}$

-Let $C_{\omega,\phi}$ be a constant such that some set $\{D_{4^{\phi-j}\omega, 2^i} : j \in \mathbb{N}_\phi\}$ of depth-robust graphs exists $\forall i \geq C_{\omega,\phi}$.

-Let C_ϕ be a constant such that $\forall i \geq C_\phi$ we have $\frac{\log(i)}{i} \geq \frac{2^{-\phi}}{3}$.

If both constants exist then $\forall i \geq \max\{C_{\omega,\phi}, C_\phi, 11\}$ the DAG $G_{\omega,\phi,i}$ has a single source and sink, size n_i and:

1. $n_i = (\phi + 1)(2^i)$.

2. $\text{indeg}(G_{\omega,\phi,i}) \leq 3i^3 + 1$.

3. $\text{cc}(G_{\omega,\phi,i}) \geq \frac{\omega^2}{144} i^3 (2^i)^{2-2^{-\phi}}$.

The proof of Lemma 6 is to use index-heavy group to solve. So based on the Lemma 6, we can know that the model's complexity in the article[?]. According to conclusion 1, we can get

$$O(i) = O(\log n) - O(\log(\phi + 1)) \quad (31)$$

We put it in the conclusion 3, we can get

$$\text{cc}(G_{\omega,\phi,i}) = \frac{\omega^2}{144} O((\log n - \log(\phi + 1))^3 (\frac{n}{\phi + 1})^{2-2^{-\phi}}) \approx n^2 (\log n)^3 \quad (32)$$

It is because that fix any $\phi \in \mathbb{N}$, so we can choose ϕ is too large to make $2 - 2^{-\phi} \rightarrow 2$, so we can get the max lower boundary is $O(n^2 (\log n)^3)$. Next, we will give our analysis about the index:cc, and we will proof cc's complexity in our model is better than its, so it can illustrate our plan is better. And here, we fixed i and ϕ to define the following three constants:

$\alpha := \lfloor (1 - 2^{-\phi})i \rfloor$ $\beta := \lfloor (1 - 2^{-\phi})i + 3\log i \rfloor$ $R := 2^{\beta-\alpha}$ // Based on these index, we can get:

$$\begin{aligned} \text{cc}(G_{\omega,\phi,i}) &= \min\{p' - \text{cost}(P) : P \in \pi\} \\ &= \min\left\{\sum_{i=0}^t |P_i| + (c-1) \sum_{i=0}^s |Q_i| : P \in \pi\right\} \\ &\geq \frac{2^{i-\alpha}}{3} (1 - 3\omega) 2^\alpha cR \\ &\approx \frac{1 - 3\omega}{3} cR 2^i \\ &\approx \frac{1 - 3\omega}{3} c 2^{\beta-\alpha} n \\ &\approx \frac{1 - 3\omega}{3} c 2^{3\log i} n \\ &\approx \frac{1 - 3\omega}{3} c i^3 n \\ &\approx O(n (\log n)^3) \end{aligned} \quad (33)$$

We can get the better complexity than original paper's model, our better complexity is $O(n (\log n)^3)$

References

- [1] Henry Corrigan-Gibbs, Dan Boneh, and Stuart E Schechter. Balloon hashing: Provably space-hard hash functions with data-independent access patterns. *IACR Cryptology ePrint Archive*, 2016:27, 2016.