

Revisiting Collision and Local Opening Analysis of ABR Hash

Chandranan Dhar ✉

Indian Statistical Institute, Kolkata, India

Yevgeniy Dodis ✉

New York University, USA

Mridul Nandi ✉

Indian Statistical Institute, Kolkata, India

Abstract

The question of building the most efficient tn -to- n -bit collision-resistant hash function H from a smaller (say, $2n$ -to- n -bit) compression function f is one of the fundamental questions in symmetric key cryptography. This question has a rich history, and was open for general t , until a recent breakthrough paper by Andreeva, Bhattacharyya and Roy at Eurocrypt'21, who designed an elegant mode (which we call ABR) achieving roughly $2t/3$ calls to f , which matches the famous Stam's bound from CRYPTO'08. Unfortunately, we have found serious issues in the claims made by the authors. These issues appear quite significant, and range from verifiably false statements to noticeable gaps in the proofs (e.g., omissions of important cases and unjustified bounds). We were unable to patch up the current proof provided by the authors. Instead, we prove from scratch the security of the ABR construction for the first non-trivial case $t = 11$ (ABR mode of height 3), which was incorrectly handled by the authors. In particular, our result matches Stam's bound for $t = 11$. While the general case is still open, we hope our techniques will prove useful to finally settle the question of the optimal efficiency of hash functions.

2012 ACM Subject Classification Security and privacy → Cryptography

Keywords and phrases ABR hash, collision resistance, local opening

Digital Object Identifier 10.4230/LIPIcs.ITC.2022.11

Funding *Yevgeniy Dodis*: Partially supported by gifts from VMware Labs and Google, and NSF grants 1815546 and 2055578.

Mridul Nandi: Partially supported by the project "Study and Analysis of IoT Security" under Government of India at R.C.Bose Centre for Cryptology and Security, Indian Statistical Institute, Kolkata

1 Introduction

The *Merkle-Damgård construction* [3, 7] is a sequential construction which is used in MD5, SHA-1 and SHA-2 and many other hash functions. On the other hand, the *Merkle tree* [6] is a parallel construction that is used in hash-based signatures (of interest due to their post-quantum security), version control systems such as git, and cryptocurrencies such as Ethereum. It is well known that the Merkle-Damgård construction and the Merkle tree are collision-resistant provided so are the compression functions. The number of compression function calls is (essentially) the same for both constructions. When we use $2n$ -to- n -bit compression functions, we can process t blocks of messages by making t or $(t - 1)$ calls to the compression function.

Although both of these widely used constructions are rather efficient, and only rely on the collision-resistance of the compression function, practical compression functions are believed to have more properties than mere collision resistance. As such, it is interesting to study the question of designing the most efficient way to build a t -to-1 collision-resistant hash function,



© Chandranan Dhar, Yevgeniy Dodis and Mridul Nandi;
licensed under Creative Commons License CC-BY 4.0
3rd Conference on Information-Theoretic Cryptography (ITC 2022).

Editor: Dana Dachman-Soled; Article No. 11; pp. 11:1–11:22



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

11:2 Revisiting Collision and Local Opening Analysis of ABR Hash

45 even if modeling the compression function as ideal (i.e. a random oracle). In particular, to see
46 whether the classical Merkle-Damgård and Merkle tree constructions can be improved under
47 such idealized modeling. This question has received a lot of attention from the cryptography
48 community, which we survey below.

49 LOWER BOUND ON THE NUMBER OF CALLS. We start with lower bounds (i.e., attacks).
50 In [2], Black et al. formally analyze the security-efficiency trade-off of compression functions,
51 showing that a $2n$ -to- n -bit compression function making a single call to a fixed-key n -bit block
52 cipher can not achieve collision resistance. Later Rogaway and Steinberger [9] generalized
53 the result for permutation-based hash. For a general hash function based on a compression
54 function, Stam [11] conjectures a lower bound on the number of compression function calls.
55 In particular, a collision with at most $2^{n(\lambda-(t-0.5)/r)}$ queries on a t -to-1 block hash function
56 can be found after making r calls to λ -to-1 block compression functions. Equivalently, for
57 optimal birthday security, the number of hash calls must be at least $r \geq (2t-1)/(2\lambda-1)$.
58 This bound is popularly known as the Stam’s bound. Stam has shown the bound for some
59 cases under a uniformity assumption. Later by Steinberger [12] and by Steinberger, Sun and
60 Yang [13], a formal proof of the Stam’s bound is shown.

61 Hence, for the most widely studied case of $\lambda = 2$, we have a lower bound $r \geq (2t-1)/3$,
62 leaving a factor 1.5 efficiency gap when compared to the Merkle-Damgård and Merkle trees.

63 UPPER BOUND ON THE NUMBER OF CALLS. For the upper bounds, much of earlier work
64 concentrated on the setting of the “non-compressing” case of $\lambda = 1$, and often focused on
65 the case of small t (e.g., $t = 2$), implicitly suggesting that — once the 2-to-1 function is
66 built, — one should do further extensions with either Merkle-Damgård and Merkle trees.
67 For example, Shrimpton and Stam [10] proposed a 2-to-1 compression function based on
68 three calls of non-compressing function, which matches Stam’s bound for $\lambda = 1$ and $t = 2$.
69 Rogaway and Steinberger [8] designed similar results when the non-compressing primitive is
70 an invertible permutation, which they also showed is optimal for this setting [9].

71 For general (large) t , Mennink and Preneel [5] also considered the non-compressing
72 case $\lambda = 1$ and proposed an elegant tree-based mode of operation making $(2t-1)$ calls to
73 the non-compressing round function, which matches Stam’s bound. Unfortunately, they
74 could only prove below-birthday security of $2^{n/3}$ queries for this construction. They also
75 conjectured that the construction achieves optimal birthday security $2^{n/2}$, but could only
76 prove it for a very restricted special-case attacker. These attacks make all their random
77 oracle calls “layer-by-layer” (as opposed to in any order). As acknowledged by the authors,
78 the simplifying assumption significantly helps with the proof of this special case and appears
79 to be with a great loss of generality. In fact, they presented evidence that their existing
80 analysis is unlikely to work for proving optimal security against unrestricted attackers.

81 Recently, two papers have appeared to tackle the compressing case $\lambda = 2$. In [4], Dodis
82 et al. optimally settled the case $t = 5$, by introducing the $T5$ construction that processes
83 five n -bit message blocks using three $2n$ -to- n -bit compression function calls, which matches
84 Stam’s bound for $t = 5$ and $\lambda = 2$. Further, they suggested extending the $T5$ construction
85 to a larger value of t using either Merkle-Damgård or Merkle trees. In both cases, they
86 already achieve non-trivial saving compared to the earlier efficiency of these modes (equal to
87 t compression calls): both variants now make roughly $3t/4$ calls to the compression functions.
88 Still, once $t > 5$, this does not match the current lower bound of $2t/3$ calls. [4] also mentioned
89 a natural, but more aggressive, variant of this extended construction for the case of Merkle
90 trees. However, they remark that this construction — even if proven collision-resistant
91 (which is open), — would lose the efficient “local opening” properties of their simpler tree

92 construction with $3t/4$ compression calls. Namely, one can no longer open one message
 93 block by only opening $O(\log t)$ internal values in the tree (as any such opening cannot have
 94 birthday security, despite satisfying correctness).

95 Finally, a breakthrough result of Andreeva, Bhattacharyya and Roy at Eurocrypt'21 [1]
 96 have claimed to settle the general case in the affirmative. They proposed a hash function ABR_l
 97 based on a perfect binary tree of height l . The hash ABR_l can process $t = (2^l + 2^{l-1} - 1)$ blocks
 98 with $r = (2^l - 1)$ calls of compression functions. This matches Stam's bound $r \geq (2t - 1)/3$.
 99 Somewhat interestingly, the ABR construction looks very similar to the tree construction of
 100 Mennink and Preneel [5] from non-compressing primitives, except all the compression calls
 101 at the leaf level now have an extra input (due to $\lambda = 2$ instead of $\lambda = 1$), while the internal
 102 calls to the compression function can also process an extra input, but using a slightly trickier
 103 rule involving two simple XOR operations. So, at least in the intuitive sense, the authors
 104 must have resolved the difficulty of [5] of dealing with general adversaries, for a construction
 105 very similar to the one of [5].

106 As an additional bonus feature, the work of [1] even claimed that the ABR_l mode also has
 107 attractive local opening properties, at the expense of slightly longer proof length ($2l$ instead
 108 of l of Merkle trees), but still having only l compression calls to verify such local opening.

109 ARE WE DONE? Unfortunately, we have found serious issues in many claims made by the
 110 authors of [1], whom we call ABR hereafter. These issues appear quite significant, and range
 111 from verifiably false statements to noticeable gaps in the proof (e.g., omissions of important
 112 cases and unjustified bounds). Unfortunately, at this stage, we are unable to fix these issues
 113 in any simple way.

114 1.1 Our Results

115 Our results can be roughly divided into 3 categories:

- 116 (1) explicit refutation of some claims made by [1];
- 117 (2) serious technical issues in the proof provided by [1];
- 118 (3) a correct (but very different from [1]) proof for the for the ABR_3 construction (i.e. $t = 11$
 119 and $r = 7$), which is incorrectly handled by ABR.

120 We detail these below.

121 LOCAL OPENING INSECURITY OF ABR. As we mentioned, ABR proposed a very efficient
 122 local opening for ABR_l . It opens about $2l$ blocks and makes l calls to verify. However, we have
 123 shown that a collision pair of the verification function can be found in $O(2^{n/2l})$ queries, which
 124 is significantly below birthday security already for $l = 2$. Hence, the suggested local opening
 125 can be broken in the above complexity. Moreover, we have shown that *any* non-trivial local
 126 opening of ABR_l satisfying a “by-pass verification” property (which is a natural class of
 127 openings that seems to include any natural opening one can think of) is broken below the
 128 birthday bound. For example, even opening $(t - 1)$ out of t inputs cannot be birthday-secure,
 129 where $t = 2^l + 2^{l-1} - 1 = 2^{\Omega(l)}$. In contrast, previous tree-like constructions (e.g., [4]) achieve
 130 birthday security with logarithmic opening length $O(l)$. This is discussed in Section 4.

131 There are two surprising aspects to this mistake. First, our attack is completely standard
 132 (using standard generalized birthday attack [14]). Second, the local opening subsection in the
 133 ABR-paper does not even mention anything about security, only focusing on the correctness
 134 of the opening. We found this quite surprising.

135 MISTAKES WITH THE MAIN PROOF. While the local opening mistake above is indisputable,
 136 the technical mistakes in the main collision resistance proof of ABR are harder to explain in

11:4 Revisiting Collision and Local Opening Analysis of ABR Hash

137 detail (at least in the Introduction, before the technical notation is developed). They are also
138 harder to state with conviction, since they often do a combination of the following pitfalls:

- 139 (a) involve imprecise statements,
 - 140 (b) state a bound which might be true, but which appears completely non-obvious to us (to
141 the extent of being the most difficult part of the proof);
 - 142 (c) point to an “analogous” earlier case, but we fail to see why the previous argument
143 generalizes;
 - 144 (d) state some bound which appears to be correct only if one makes some restricting assumption
145 on the attacker (but no such assumptions are made by the authors, who claim a fully
146 general result!);
 - 147 (e) silently omitting an important special case of the proof (i.e., the proof is non-exhaustive).
- 148 The totality of these issues make the proof presented by [1] at best unverifiable, and at worst
149 incorrect. In particular, we still believe that the end result is correct, but fixing it would
150 require a substantially harder proof.

151 At a very high level, the correct collision analysis for a tree-based function like ABR_l
152 is complex mostly due to the *adaptive nature* of queries, and the queries made to different
153 layers in the tree might not come in monotone order (i.e., may not be in order of the level of
154 the nodes). Indeed, this is precisely the reason why the earlier birthday security result of
155 Mennink and Preneel [5] only held for “in order” adversaries. Fortunately, the outputs of the
156 leaf nodes can be given beforehand, as the input of those has no role in finding a collision.
157 More formally, we can make a simple argument to force the attacker “evaluate” the first
158 layer compression calls before any of the subsequent calls as follows. We give the attacker q
159 random outputs (where q is the total number of queries made by the attacker) at the very
160 beginning, but allow the adversary to *arbitrarily label* the corresponding input values at any
161 point in the game. This is fine, since those input values do not participate in any other
162 computation, but now all the outputs in the first layer are known before a single compression
163 call is made to the lower layers. This allows for relatively simple analysis for the special case
164 $l = 2$, and the authors of [1] indeed start with the correct analysis of this special case.¹

165 Unfortunately, this argument completely fails after the first layer. (Indeed, handling this
166 case will be one of the most difficult parts of our analysis, when we provide a correct proof
167 for $l = 3$ in this paper.) In particular, we see the following high-level issues with the proof
168 presented by [1] for $l \geq 3$. (More lower-level issues are discussed in Section 5.3 in the paper.)

- 169 1. ABR claimed a relation between collision and the number of computable hash outputs
170 (termed as load). We will show in Section 5.4 that the relation is not true in general
171 by giving a counterexample. This seems to hold for ABR if queries to the root node are
172 performed at the end (which is the case for ABR_2). However, it seems non-obvious to us
173 why a similar relation holds when the adversary makes out-of-order queries.
- 174 2. We have also found issues while bounding load. ABR consider “input multi-collision”
175 for every node up to $O(n)$. However, due to the multiplicative nature of the number of
176 multi-collisions as one goes down in the tree, we find that $O(n^i)$ multi-collision must
177 be considered for the nodes at the i -th level. This would degrade the bound for load
178 claimed by ABR, and invalidate the claimed birthday security at the end (unless the

¹ Another correct proof for $t = 5$ (corresponding to tree depth $l = 2$) was made for the $T5$ compression function by [4]. Interestingly, the authors did not notice the simplifying non-adaptivity argument above, and had to work relatively hard to handle out-of-order queries (e.g., it involved a careful expectation analysis and applying Markov’s inequality; see proof of Proposition 5 in [4], which is over a page). This shows that handling out-of-order attackers is indeed highly non-trivial.

179 number of levels i is constant, in which case one can hide the extra n^i bound in the
 180 “ O -tilde”-notation). This will be discussed in Step 1 of Section 5.3.

181 **3.** In fact, even if the load analysis is somehow fixed, ABR seem to consider the last query
 182 happens in the final node (or at the node where the load is considered). This is effectively
 183 equivalent to in order adversaries, but does not seem to be the case for general attackers.
 184 See Step 2 of Section 5.3.

185 **4.** Moreover, both messages of a collision pair can be generated due to a single query response
 186 (termed as *twin collision pair*). ABR completely ignore this case. This is discussed in
 187 detail in the last paragraph of Section 5.3.

188 We leave a more detailed explanation of these (and other issues (a)-(e)) later in the paper.

189 COLLISION ANALYSIS OF ABR₃. On a positive, our main technical result shows that the
 190 ABR₃ construction for $t = 11$ indeed achieves birthday security (roughly $n^5 q^2 / 2^n$, where q
 191 is the number of compression function queries) with an optimally small number of $r = 7$
 192 compression calls (see Section 6). While forming only the first step in recouping the incorrect
 193 results of [1], we are optimistic that our approach could be extended to finally settle the
 194 general case correctly. For example, compared to best known correct proofs for $t = 5$ (e.g.,
 195 ABR₂ from [1], or the $T5$ compression function from [4]), we can no longer assume that the
 196 second layer calls to the compression function are made before all the third-layer calls, which
 197 is the main (unresolved) difficulty in the work of [5], and one of the key mistakes in the
 198 analysis of [1] (as we explained above). Thus, our proof is the first which handles non-trivial
 199 “out-of-order” adversaries correctly.

200 We also hope our proof of ABR₃ provides a sharp contrast to the flawed proof of [1], even
 201 for this special case. For example, we already mentioned handling general “out-of-order”
 202 adversaries. In a different vein, we also consider the twin-collision analysis for ABR₃ which is
 203 completely missing from [1]. This analysis requires a non-trivial multi-collision analysis on a
 204 sum of our compression functions, and we also need to bound some other failure events to
 205 analyze the non-twin collision security of ABR₃. None of these arguments appeared in [1].

206 **2 Security Definitions**

207 **2.1 Notations**

208 We call elements of $\{0, 1\}^n$ blocks. A k -to- r (block) function or random oracle has domain
 209 $\{0, 1\}^{kn}$ and range $\{0, 1\}^{rn}$. We write the set $[k] = \{1, 2, \dots, k\}$. A partial function τ from D
 210 to R is a subset $\tau \subseteq D \times R$ such that for every $x \in D$, there are at most one y with $(x, y) \in \tau$.
 211 We define domain $\text{dom}(\tau) := \{x : \exists y, (x, y) \in \tau\}$ and range $\text{ran}(\tau) = \{y : \exists x, (x, y) \in \tau\}$ of a
 212 partial function τ . We use the shorthand notation $A \cup x$ and $A \setminus x$ to denote $A \cup \{x\}$ and
 213 $A \setminus \{x\}$ respectively. For any q -tuple x^q , we define $\text{mc}(x^q) = \max_a |\{i : x_i = a\}|$. For two
 214 lists \mathcal{L}_1 and \mathcal{L}_2 , we define $\text{mc}(\mathcal{L}_1 \oplus \mathcal{L}_2) = \max_a |\{(i, i') : L_i \oplus L_{i'} = a, L_i \in \mathcal{L}_1, L_{i'} \in \mathcal{L}_2\}|$. It
 215 can be similarly extended for xor of more than two lists.

216 **2.2 Generic Hash Mode**

217 Let H^f be a t -to-1 hash function which uses an n -bit compression function (i.e. λ -to-1
 218 compression function f for some $\lambda > 1$) as an oracle. Note that a mode can use more than
 219 one compression functions f_1, \dots, f_r . However, as we analyze in the random oracle model,
 220 independent random oracles can be obtained from a single random oracle with a little bit
 221 larger domain by using the standard domain separation method. In this paper, we only
 222 consider fixed-length input and also assume r is the same for all messages. Moreover, the

11:6 Revisiting Collision and Local Opening Analysis of ABR Hash

hash function calls f_i on i -th call and so the domains of every call are separated by domain separation. We also denote the family $f := (f_i : i \in [r])$ by f and we call λ -to-1 r r.o. (random oracle). We denote $\tau_{\mathbb{H}}(M \mid f) := \{((1, x_1), y_1), \dots, ((r, x_r), y_r)\}$ where x_i denotes the input of i -th call of its oracle tuple while computing $\mathbb{H}^f(M)$ and $y_i = f_i(x_i) := f(i, x_i)$. A λ -to-1 **transcript** τ is a partial function from $[r] \times \{0, 1\}^{\lambda n}$ to $\{0, 1\}^n$. For a λ -to-1 r r.o. f , we have

$$\forall (i, x) \notin \text{dom}(\tau), y \in \{0, 1\}^n, \text{Prob}(f(i, x) = y \mid \tau \subseteq f) = 2^{-n}.$$

► **Definition 1** (transcript-based hash computation). *Given a partial function $\tau \subseteq f$, let $\mathbb{H}^\tau = \{(M, \mathbb{H}^f(M)) : \tau_{\mathbb{H}}(M \mid f) \subseteq \tau\}$ be a partial hash function. In other words, \mathbb{H}^τ consists of all pairs (M, z) such that $\mathbb{H}^f(M)$ can be computed by simply using the transcript τ and z is the final value. The elements of the set $\text{dom}(\mathbb{H}^\tau)$ are called τ -computable messages. As $\tau \subseteq f$, we must have $\mathbb{H}^\tau \subseteq \mathbb{H}^f$.*

2.3 Collision Game

Let \mathcal{A} be an adversary having oracle access of f which makes q queries to each f_i adaptively. As we assume an unbounded time adversary, there is no loss in assuming that \mathcal{A} is deterministic. Thus, the i -th query (x_i, v_i) of \mathcal{A} depends on τ^{i-1} (the transcript of query-responses after $(i-1)$ queries). After the query-response phase, \mathcal{A} returns a pair of distinct messages (M, M') such that both M, M' are transcript-computable. We say $\text{coll}_{\mathbb{H}}$ holds if $\mathbb{H}^\tau(M) = \mathbb{H}^\tau(M')$, called a *computable collision pair*. We define $\text{Adv}_{\mathbb{H}^f}^{\text{coll}}(\mathcal{A}) := \Pr(\text{coll}_{\mathbb{H}})$.

► **Definition 2** (cross collision). *Let \mathbb{H} and \mathbb{H}' be two hash functions. A cross-collision τ -computable pair is a pair (M, M') (not necessarily distinct) such that $\mathbb{H}^\tau(M) = \mathbb{H}'^\tau(M')$. We denote $\text{coll}_{\mathbb{H}, \mathbb{H}'}^\tau := \{M \in \text{dom}(\mathbb{H}^\tau) : \exists M', \mathbb{H}^\tau(M) = \mathbb{H}'^\tau(M')\}$.*

2.4 Local Opening

We now define the local opening security of a hash function output (viewed as a commitment of a message). Given a hash function mode \mathbb{H}^f , a local opening Open^f for \mathbb{H} maps a pair (M, i) to $\pi = (m_i, i, \pi')$ (called proof) where $M = (m_1, m_2, \dots, m_c)$ is a message (a tuple of blocks) and $i \in [c]$ is an index.

CORRECTNESS OF LOCAL OPENING. There is an efficient function Ver^f such that for all message M , all index i , $\text{Ver}^f(\text{Open}^f(M, i), \mathbb{H}^f(M)) = 1$.

SECURITY OF LOCAL OPENING. In the local opening security, the adversary wins if it produces an output h corresponding to two contradicting local openings for some position i .

► **Definition 3** (local opening advantage). *Let \mathbb{H} be a hash function and Open be a correct local opening for \mathbb{H} with verification function Ver . For any adversary \mathcal{A} , we define the local opening advantage as*

$$\text{Adv}_{\mathbb{H}}^{\text{local}}(\mathcal{A}) = \Pr \left[\text{Ver}(i, m, \pi, h) = \text{Ver}(i, m', \pi', h) = 1, m \neq m' \right. \\ \left. \mid (i, m, m', \pi, \pi', h) \leftarrow \mathcal{A}^f \right]$$

BY-PASS HASH COMPUTATION. We say that \mathbb{H} has a *by-pass computation* $(\mathbb{H}_i : i \in [c])$ corresponding to a local opening Open if for all $M, i \in [c]$,

$$\mathbb{H}_i^f(\text{Open}^f(M, i)) = \mathbb{H}^f(M).$$

263 In other words, given a proof (output of the `Open`) and the message block for the index (for
 264 which the proof is produced), we can compute the hash output of the message (without
 265 knowing the other blocks of the message). The verification algorithm simply checks whether
 266 the hash value computed through the by-pass hash is the same as what was committed
 267 before. As f is treated as an oracle, it is natural to assume that for all M and for all i ,

$$268 \quad \tau_{\text{Open}}(M, i \mid f) \cup \tau_{\text{H}_i}(\text{Open}^f(M, i) \mid f) = \tau_{\text{H}}(M \mid f).$$

269 We now define the *inter-collision* advantage for by-pass computation H_i as

$$270 \quad \text{Adv}_{\text{H}_i}^{\text{coll}^*}(\mathcal{A}) = \Pr \left[\text{H}_i(m, \pi) = \text{H}_i(m', \pi') \text{ and } m \neq m' \mid (m, \pi, m', \pi) \leftarrow \mathcal{A}^f \right].$$

272 Thus, it is the same as the collision game, except that the adversary needs to find a collision
 273 pair for which $m \neq m'$. Suppose \mathcal{A} finds a collision pair $((m, \pi), (m', \pi'))$ for H_i , and let
 274 $h = \text{H}_i(m, \pi)$. Then \mathcal{A} can commit h and later on, it can successfully open for either of two
 275 messages m and m' as required. Now we make the following simple observation

$$276 \quad \text{Adv}_{\text{H}}^{\text{local}}(q') = \max_{\mathcal{A}} \max_i \text{Adv}_{\text{H}_i}^{\text{coll}^*}(\mathcal{A}). \quad (1)$$

278 The above observation (see [4] for details) helps us to reduce the local opening security to
 279 inter-collision security problem for the by-pass hash family.

280 2.5 Stam's Tradeoff between Security and Performance

281 Stam's bound states that there always exists a collision attack with at most $2^{n(\lambda - (t-0.5)/r)}$
 282 queries on a t -to-1 block hash function making r calls to λ -to-1 block compression functions.

283 3 Re-introduction of the ABR Hash due to [1]

284 We first start by defining a generalized tree hash structure, and then re-introduce the ABR
 285 Hash as a special tree hash, as opposed to introducing as it is in [1]. This is because we feel
 286 some things have not been properly defined by the authors there, and these issues need to
 287 be addressed properly.

288 A *full binary tree* (FBT) is a binary tree in which every node v other than the leaves has
 289 two children, denoted as v_L (left child) and v_R (right child). A *perfect binary tree* (PBT) is a
 290 full binary tree in which all the leaf nodes are at the same level (called height of the tree).

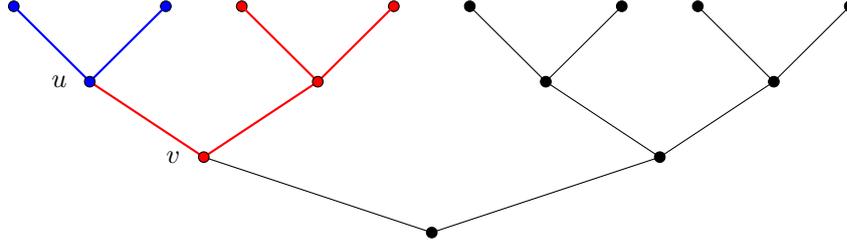
291 ► **Example 4** (perfect binary tree of height l). Let l be a fixed positive integer and \mathcal{T} be
 292 a perfect binary tree of height l over all vertices (j, b) , $j \in [l], b \in [2^{l-j}]$ with $(l, 1)$ being
 293 the root. For every two vertices (j, b) and $(j+1, \lceil b/2 \rceil)$, we associate an edge. We call
 294 $(j-1, 2b-1)$ and $(j-1, 2b)$ the left and right child of (j, b) respectively. Note that $\mathcal{T} = \mathcal{T}_{(l,1)}$.

295 3.1 Some Notations and Definitions on Binary Trees

296 For a binary tree \mathcal{F} , let $\mathcal{L}_{\mathcal{F}}$ and $V(\mathcal{F})$ denote the set of leaf nodes and all nodes of \mathcal{F}
 297 respectively. Any non-leaf node is called an intermediate node. For a non-root intermediate
 298 node v of \mathcal{F} , we consider the following two full binary trees:

- 299 1. \mathcal{F}_v : the full binary sub-tree rooted at v .
- 300 2. \mathcal{F}_{-v} : the sub-tree $(\mathcal{F} \setminus \mathcal{F}_v) \cup v$.

11:8 Revisiting Collision and Local Opening Analysis of ABR Hash



■ **Figure 1** In this figure, \mathcal{F}_v is the sub-tree rooted at v , i.e. the union of the red and blue sub-trees, \mathcal{F}_{-v} is the black sub-tree, and \mathcal{F}_{v-u} is the red sub-tree.

301 For a tree \mathcal{F} , and a vertex v of \mathcal{F} , we write V_v , \mathcal{L}_v and V_v^* to denote the set of all nodes, leaf
 302 nodes and intermediate (non-leaf) nodes respectively for the tree \mathcal{F}_v . For any $u \in V_v^* \setminus v$, we
 303 write $\mathcal{F}_{v-u} = (\mathcal{F}_v \setminus \mathcal{F}_u) \cup u$. We write V_{v-u} to denote the set of vertices of \mathcal{F}_{v-u} . For the
 304 sake of notational simplicity we ignore the suffix v when v is the root. In this section we only
 305 consider trees of the form \mathcal{F}_v and \mathcal{F}_{v-u} . Refer to Figure 1 for a pictorial representation.

306 To each node $v \in V$ of a perfect binary tree \mathcal{T} , an independent 2-to-1 block compression
 307 function (modeled as a random oracle) f_v is assigned. We use the notation f to denote the
 308 collection of random oracles $\{f_v : v \in \mathcal{T}\}$.

309 ► **Definition 5** (message for tree hash). *A message m for any full binary sub-tree \mathcal{F} of a
 310 perfect binary tree \mathcal{T} having the same root is a function $m : V(\mathcal{F}) \rightarrow \{0, 1\}^n \cup \{0, 1\}^{2n}$ such
 311 that for all $u \in \mathcal{L}_{\mathcal{F}} \cap \mathcal{L}_{\mathcal{T}}$, $m(u) \in \{0, 1\}^{2n}$, otherwise, $m(u) \in \{0, 1\}^n$. A complete message
 312 m is a message at the root of \mathcal{T} .*

313 Thus, for every leaf node of \mathcal{F} (which is also a leaf node of the perfect binary tree), we
 314 associate $2n$ bit messages. For all other vertices, we associate an n bit message. We write
 315 $\mathbb{M}_{\mathcal{F}}$ to denote the set of all messages for \mathcal{F} . We simply write \mathbb{M}_v and \mathbb{M}_{v-u} instead of $\mathbb{M}_{\mathcal{T}_v}$
 316 and $\mathbb{M}_{\mathcal{T}_{v-u}}$ respectively.

317 For a message m for \mathcal{T}_v (also called m at the node v), and $u \in V_v$, we write $m|_u = m|_{\mathcal{T}_u}$,
 318 the message restricted to \mathcal{T}_u . Similarly, we write $m_L := m|_{v_L}$ and $m_R := m|_{v_R}$. We also
 319 write $m|_{v-u \rightarrow h}$ to denote a message for \mathcal{T}_{v-u} which is same as the restricted function $m|_{\mathcal{T}_{v-u}}$,
 320 except at u , where it assigns h (instead of $m(u)$). In the context of our work, this basically
 321 means we replace the message $m(u)$ at node u by the intermediate hash output of \mathcal{T}_u , the
 322 tree rooted at u , and consider the message for the remaining tree, \mathcal{T}_{v-u} .

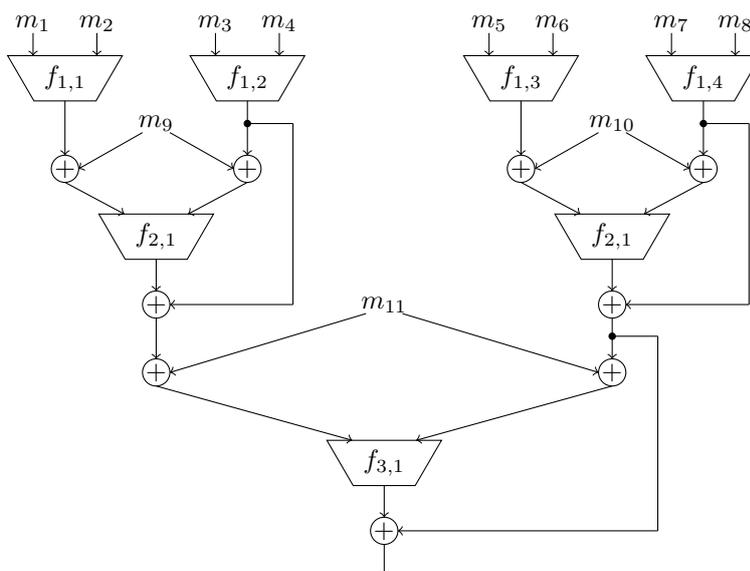
323 ► **Definition 6** (Generalized Tree Hash). *Let \mathcal{F} be a full binary sub-tree of a perfect binary
 324 tree \mathcal{T} and let $m \in \mathbb{M}_{\mathcal{F}}$. For every $v \in \mathcal{F}$, we associate an intermediate hash output O_v and
 325 an intermediate input I_v recursively as follows:*

- 326 1. $v \in \mathcal{L}_{\mathcal{F}} \setminus \mathcal{L}_{\mathcal{T}}$, $|m(v)| = n$: $O_v = m(v)$ and there is no input,
- 327 2. $v \in \mathcal{L}_{\mathcal{F}} \cap \mathcal{L}_{\mathcal{T}}$, $|m(v)| = 2n$: $O_v = f_v(m(v))$, $I_v = m(v)$,
- 328 3. otherwise: $|m(v)| = n$ and we define

$$329 \quad I_v = (O_{v_L} \oplus m(v), O_{v_R} \oplus m(v)), \text{ and } O_v = f_v(O_{v_L} \oplus m(v), O_{v_R} \oplus m(v)) \oplus O_{v_R}.$$

330 O_ω is the final hash output corresponding to \mathcal{F} where ω is the root of \mathcal{F} . We also call I_ω
 331 final input.

332 Let us see what this means. If $\mathcal{F} = \mathcal{T}$, the above definition implies that for a leaf node v ,
 333 the message at v , which itself is the input, is $2n$ bits long, and the output is just $f_v(m(v))$,
 334 where f_v is the 2-to-1 block compression function attached to it, and for an intermediate



■ **Figure 2** ABR of height 3

node, the message is n bits long, and the input and output are as defined above. If \mathcal{F} is a proper sub-tree of \mathcal{T} , then there might exist vertices, which are leaves of \mathcal{F} , but not of \mathcal{T} . For such a vertex v , the message is n bits long, and the message itself is considered the output of the vertex. This vertex doesn't have any input.

THE ABR HASH FUNCTION. The ABR hash is the hash output based on a perfect binary tree \mathcal{T} of height l . In terms of Definition 6, the case $\mathcal{F} = \mathcal{T}$ corresponds to a ABR tree, and the final hash output is the ABR hash. Thus, ABR_l hash is a $(2^l + 2^{l-1} - 1)$ -to-1 block hash function, $l > 1$. We refer to Figure 2 for a pictorial view of ABR with $l = 3$. For a trivial tree $\mathcal{F} = \{w\}$, with a message $m(\omega) \in \{0, 1\}^{2^n}$, $\mathcal{F}(m) = f_\omega(m)$.

We write $\text{H}^\tau(m)$ and $\text{in}^\tau(m)$ to denote the transcript based hash and the final input respectively, whenever defined for the message m for a tree \mathcal{F} . If $\text{H}^\tau(m)$ is defined we call m τ -computable or simply computable message. We write \perp to mean that it is undefined. Note that a tree is uniquely determined from the message. We write dom_v^τ and dom_{v-u}^τ to denote the set of all computable messages at v and for \mathcal{T}_{v-u} respectively. Similarly, we write ran_v^τ and ran_{v-u}^τ to denote the set of all computable hashes at v and for \mathcal{T}_{v-u} respectively. The size of the set ran_v^τ , called *load at v* , is denoted as $L_{\mathcal{T},v}$.

4 Local Opening Analysis of ABR Hash Function

In section 3, we have defined hash function based on a tree \mathcal{F} for a message over the tree \mathcal{F} . In this section, we consider a variant of the message function and a hash function for the variant message. This is required to properly define the local opening of the ABR tree.

MESSAGE FOR A FULL BINARY TREE. Let \mathcal{F} be a full binary tree and $L \subseteq \mathcal{L}_{\mathcal{F}}$. Let $\mathcal{M}_{\mathcal{F},L}$ be the set of all functions $m : V(\mathcal{F}) \rightarrow \{0, 1\}^n \cup \{0, 1\}^{2^n}$ such that for all $v \in \mathcal{L}_{\mathcal{F}} \setminus L$, $m(v) \in \{0, 1\}^{2^n}$ and for all other vertices v , $m(v) \in \{0, 1\}^n$. We call m a message (or a message function) for \mathcal{F} .

► **Definition 7** (Generalized Tree Hash, a variant). Let $m \in \mathcal{M}_{L,\mathcal{F}}$ be a message function for \mathcal{F} . For every $v \in \mathcal{F}$, the intermediate hash output O_v is defined recursively as follows:

11:10 Revisiting Collision and Local Opening Analysis of ABR Hash

- 361 1. $v \in L$, $|m(v)| = n$: $O_v = m(v)$,
- 362 2. $v \in \mathcal{L}_{\mathcal{F}} \setminus L$, $|m(v)| = 2n$: $O_v = f_v(m(v))$, $I_v = m(v)$,
- 363 3. $v \notin \mathcal{L}_{\mathcal{F}}$: we define

$$364 \quad I_v = (h_1 \oplus m(v), h_2 \oplus m(v)) \text{ and } O_v = f_v(h_1 \oplus m(v), h_2 \oplus m(v)) \oplus h_2,$$

365 where $h_1 = O_{v_L}$ and $h_2 = O_{v_R}$.

366 The hash output corresponding to \mathcal{F} is defined as $\mathcal{F}^f(m) := O_\omega$ where ω is the root of \mathcal{F} .
 367 We also call $I_\omega := \mathcal{F}_{\text{in}}^f(m)$ final input. It is clear from the definition that for any node $v \notin L$,
 368 $\mathcal{F}_v^f(m|_v) = O_v$ and $\mathcal{F}_{v,\text{in}}^f(m|_v) = I_v$.

369 Visualizing the tree is not difficult. As an example, when $\mathcal{F} = \text{ABR}_3$, we have Figure 2,
 370 where L is a subset of the leaf nodes, say $(1, 1)$ and $(1, 2)$. We now define local opening of
 371 the Generalized Tree Hash.

372 ► **Definition 8.** Let m be a message for a perfect binary tree \mathcal{T} . For any full binary sub-tree
 373 \mathcal{F} and a set $\mathcal{L}_{\mathcal{F}} \setminus \mathcal{L}_{\mathcal{T}} \subseteq L \subseteq \mathcal{L}_{\mathcal{F}}$, we define a message $m' := \text{Open}_{\mathcal{F},L}^f(m) \in \mathcal{M}_{\mathcal{F},L}$ for \mathcal{F} as
 374 follows.

- 375 1. $v \in L$: $m'(v) = \mathcal{T}_v^f(m_v)$.
- 376 2. Otherwise: $m'(v) = m(v)$.

377 Now, we first analyze the local opening security of ABR_l proposed by [1] and then show
 378 that no non-trivial opening of ABR can achieve birthday bound security.

379 4.1 Local Opening Analysis of ABR Hash due to [1]

380 We describe the by-pass hash corresponding to the message block m_1 for ABR_l . It is based
 381 on the full sub-tree \mathcal{F} consisting of nodes $\{(i, 1) : i \in [l]\} \cup \{(i, 2) : i \in [l-1]\}$ and
 382 $L = \{(1, 2), (2, 2), \dots, (l-1, 2)\}$. Refer to Figure 3. Note that the number of blocks in
 383 $\text{Open}_{\mathcal{F},L}(m)$ is $2l$, and in the sub-tree \mathcal{F} corresponding to $\text{Open}_{\mathcal{F},L}(m)$, the number of calls
 384 to underlying compression function f is l . According to Stam's bound, there exists a collision
 385 attack with at most $2^{n/2l}$ queries. We give an attack that matches this bound.

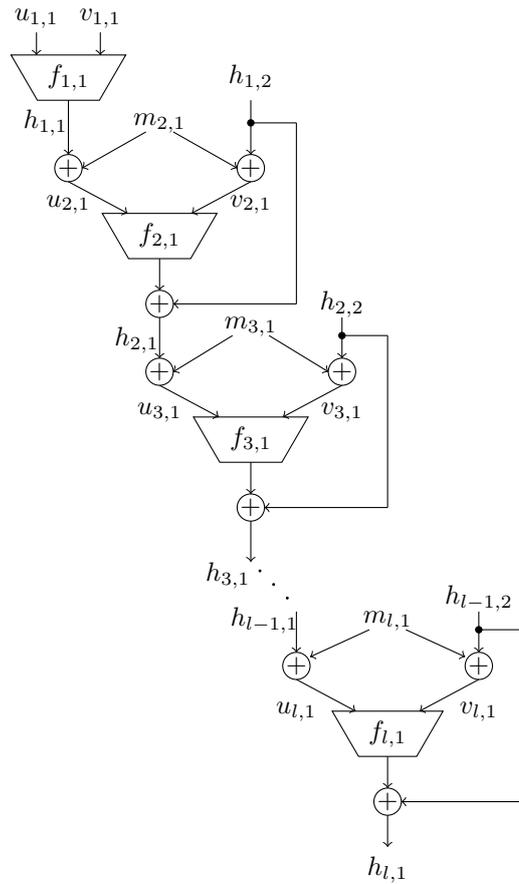
386 Let $I(i, 1) = (u_{i,1}, v_{i,1})$ be the input of $f_{i,1}$ and let $y_{i,1}$ be the output. Let $h_{i,2}$ be the
 387 message for node $(i, 2)$. Let $h_{i,1} = y_{i,1} \oplus h_{i-1,2}$ for $i > 1$ and $h_{1,1} = y_{1,1}$. Then, $h_{i,1}$ is the
 388 output at node $(i, 1)$. Also, let $m_{i,1}$ be the message associated with a non-leaf node $(i, 1)$.
 389 We wish to find a collision at the output of node $(l, 1)$, i.e. we need to find two messages m'
 390 and m'' for \mathcal{F} such that $\mathcal{F}_{(l,1)}(m') = \mathcal{F}_{(l,1)}(m'')$. Given any message for \mathcal{F} , the output at
 391 node $(l, 1)$ is given by $h_{l,1}$.

392 Note that $h_{1,1} = f_{1,1}(u_{1,1}, v_{1,1})$. After computing $h_{i-1,1}$, we proceed to compute $h_{i,1}$.
 393 We note that $h_{i-1,2}$ is a message block for \mathcal{F} . The input at node $(i, 1)$, $I(i, 1) = (h_{i-1,1} \oplus$
 394 $m_{i,1}, h_{i-1,2} \oplus m_{i,1}) = (u_{i,1}, v_{i,1})$ and the output at node $(i, 1)$ is:

$$\begin{aligned} 395 \quad h_{i,1} &= f_{i,1}(I(i, 1)) \oplus h_{i-1,2} = f_{i,1}(u_{i,1}, v_{i,1}) \oplus h_{i-1,2} \\ 396 &= f_{i,1}(u_{i,1}, v_{i,1}) \oplus u_{i,1} \oplus v_{i,1} \oplus h_{i-1,1} \\ 397 &= g_{i,1}(u_{i,1}, v_{i,1}) \oplus h_{i-1,1} \end{aligned}$$

399 where $g_{i,1}(u_{i,1}, v_{i,1}) = f_{i,1}(u_{i,1}, v_{i,1}) \oplus u_{i,1} \oplus v_{i,1}$. By induction, the final hash computation is

$$400 \quad h_{l,1} = g_{l,1}(u_{l,1}, v_{l,1}) \oplus g_{l-1,1}(u_{l-1,1}, v_{l-1,1}) \oplus \dots \oplus g_{1,1}(u_{1,1}, v_{1,1}).$$



■ **Figure 3** A specific local opening of ABR_3 .

401 Since the functions $f_{i,1}$ are random and independent so are $g_{i,1}$'s. Thus $h_{l,1}$ is the XOR of l
 402 random functions. Thus, a collision is expected at node $(l, 1)$ with $2^{n/2^l}$ queries. One can
 403 also apply a generalized birthday attack with complexity $2^{n/(1+\lceil \log 2^l \rceil)}$.

404 Now, let us look at the target collision resistance of the above local opening of ABR_l .
 405 Target Collision Resistance describes the ability of an adversary to find a second pre-image
 406 for a fixed message. Target collision resistance has many practical applications. For example,
 407 if a client sends a file F to the server and then wants the server to send part of the file F_i
 408 along with a proof of correctness then, as long as the server does not control the choice of
 409 the file F , the server would need to find a targeted collision to break security and reveal an
 410 incorrect value F'_i .

411 Here, for a fixed message m , the final hash computation $h_{l,1}$ is fixed. Hence, for target
 412 collision resistance we wish the XOR of l random functions to collide with this value of $h_{l,1}$.
 413 This collision is expected with $2^{n/l}$ queries.

414 4.2 Decomposition of ABR Hash

415 Now we decompose ABR hash computation on \mathcal{T} through a full binary proper sub-tree \mathcal{F}
 416 sharing the same root and a set L .

417 ► **Lemma 9** (decomposition lemma for any full binary tree). *For all full binary sub-tree \mathcal{F} of*

11:12 Revisiting Collision and Local Opening Analysis of ABR Hash

418 a perfect binary tree \mathcal{T} and a set of nodes $\mathcal{L}_{\mathcal{F}} \setminus \mathcal{L}_{\mathcal{T}} \subseteq L \subseteq \mathcal{L}_{\mathcal{F}}$, we have

$$419 \quad \mathcal{T}^f = \mathcal{F}^f \circ \text{Open}_{\mathcal{F},L}^f.$$

420 **Proof.** Let m be a message for \mathcal{T} . $\mathcal{T}^f(m)$ represents the hash output based on the perfect
421 binary tree \mathcal{T} . For any node v of \mathcal{T} , the restricted message over \mathcal{T}_v is m_v . Hence, $\mathcal{T}^f(m)$
422 computes $\mathcal{T}_v^f(m_v)$ for all nodes $v \in \mathcal{T}$.

423 For any full binary sub-tree \mathcal{F} of \mathcal{T} , $m' = \text{Open}_{\mathcal{F},L}^f$ is defined as above. For any $v \in L$:
424 $m'(v) = \mathcal{T}_v^f(m_v)$. We calculate the hash outputs for the restricted messages on these nodes
425 first. Since for all other $v \in \mathcal{F}$, $m'(v) = m(v)$, and \mathcal{F} is a sub-tree of \mathcal{T} , $\mathcal{F}^f(m')$ actually
426 computes $\mathcal{T}_v^f(m_v)$ for all $v \in \mathcal{F} \setminus \mathcal{L}_{\mathcal{F}}$. Thus, $\mathcal{F}^f \circ \text{Open}_{\mathcal{F},L}^f(m)$ also computes $\mathcal{T}_v^f(m_v)$ for all
427 nodes $v \in \mathcal{T}$ and produces the same output $\mathcal{T}^f(m)$. \blacktriangleleft

428 If $\mathcal{F} = \mathcal{T}$ and $L = \emptyset$ then $\text{Open}_{\mathcal{F},L}^f(m) = m$. For any other proper local opening we
429 cannot ensure birthday bound security. We prove the following theorem:

430 **► Theorem 10.** *No non-trivial opening of ABR can achieve birthday bound security.*

431 **Proof.** Stam's bound states that there exists a collision attack with at most $2^{n(\lambda-(t-0.5)/r)}$
432 queries on a t -to-1 block hash function making r calls to λ -to-1 block compression functions.
433 We have $\lambda = 2$. If we want to achieve $2^{n/2}$ collision security, $t \leq 1.5r + 0.5$. In other words,
434 if $t > 1.5r + 0.5$, then we have a collision attack with query complexity $2^{\frac{n}{2}(1-\delta/r)}$, $\delta :=$
435 $t - 1.5r - 0.5$.

436 For ABR of height l , we have $t = 2^l + 2^{l-1} - 1$ and $r = 2^l - 1$. This satisfies $t = 1.5r + 0.5$,
437 and it is optimal. We show that for any non-trivial opening $\text{Open}_{\mathcal{F},L}$ of ABR, \mathcal{F} satisfies
438 $t > 1.5r + 0.5$. Let us consider the simplest non-trivial opening, corresponding to $L = \{(1, 1)\}$.
439 Then, for $m = (m_1, m_2, m')$, where m_1, m_2 are the first two message blocks and m' is the
440 remaining part, $\text{Open}_{\mathcal{F},L}(m) = (f_{1,1}(m_1, m_2), m')$. Then, $t = 2^l + 2^{l-1} - 2$, and $r = 2^l - 2$
441 ($f_{1,1}$ is not called). This satisfies $t > 1.5r + 0.5$. If $\text{Open}_{\mathcal{F},L}$ consists of only one sub-tree
442 computation of height h , then for \mathcal{F} , we have $t = (2^l + 2^{l-1} - 1) - (2^h + 2^{h-1} - 1) + 1$ and
443 $r = 2^l - 2^h$, which satisfies $t > 1.5r + 0.5$.

444 A general opening $\text{Open}_{\mathcal{F},L}$ of ABR may consist of more than one complete sub-tree
445 computation. Let the number of complete sub-tree computations in $\text{Open}_{\mathcal{F},L}$ be k , and for
446 each $1 \leq i \leq k$, let h_i be the height of the i -th sub-tree. Then, for \mathcal{F} , we have

$$447 \quad t = (2^l + 2^{l-1} - 1) - \sum_{i=1}^k (2^{h_i} + 2^{h_i-1} - 1) + k, \quad r = (2^l - 1) - \sum_{i=1}^k (2^{h_i} - 1).$$

448 It can be easily seen that $t > 1.5r + 0.5$. Thus, no non-trivial opening of ABR can achieve
449 birthday bound security. \blacktriangleleft

5 Collision Analysis of ABR hash

450 In this section, we first define certain items which will be required to analyze the collision.

451 **► Definition 11** (input multi-collision). *For any $x \in \{0, 1\}^n$, let $\text{MC}_v^\tau(x)$, called input multi-*
452 *collision set at v (with x as input multi-collision value), denote the set of all messages m at*
453 *v with $\text{in}^\tau(m) = x$. also, let*

$$454 \quad \text{mc}_v^\tau(x) = |\text{MC}_v^\tau(x)|, \quad \text{mc}_v^\tau = \max_{x \in \{0,1\}^n} \text{mc}_v^\tau(x).$$

455 *When v is the root node, we skip the notation v .*

457 We define the newly generated messages and the hashes at a node v due to addition of the
458 query-response (x, y) to the transcript τ as

$$459 \quad \text{New}_v^\tau(x, y) := \text{dom}_v^{\tau \cup (x, y)} \setminus \text{dom}_v^\tau, \quad \text{NewH}_v^\tau(x, y) := \text{ran}_v^{\tau \cup (x, y)} \setminus \text{ran}_v^\tau.$$

460 Clearly, $\text{NewH}_v^\tau(x, y) = \text{H}^{\tau \cup (x, y)}(\text{New}_v^\tau(x, y))$ (image set of $\text{H}^{\tau \cup (x, y)}$ for the domain $\text{New}_v^\tau(x, y)$).

461 Note that x need not be queried at v . However, to have a new computable message, x should
462 be queried at some node, say u , in \mathcal{T}_v . Analyzing the behavior of the set $\text{New}_v^\tau(x, y)$ (or its
463 size) is easy when $u = v$ or when u is one of the children of v . However, it becomes more
464 complex when u is far away from v .

465 ■ Case $u = v$: $\text{New}_v^\tau(x, y) = \text{MC}_v^\tau(x)$ (and does not depend on y) and we call these messages
466 freshly generated **immediate** messages.

467 ■ Case $u \in \mathcal{T}_v \setminus v$: The newly generated messages at v is

$$468 \quad \text{New}_v^\tau(x, y) = \{m|_v : \text{in}^\tau(m|_u) = x, m|_{v-u \rightarrow h} \in \text{dom}_{v-u}^\tau, h = y \oplus \text{H}^\tau(m|_{u_R})\}.$$

469 So, we have $\mathbb{E}_y(|\text{New}_v^\tau(x, y)|) = \frac{\text{mc}_u^\tau(x) \times |\text{dom}_{v-u}^\tau|}{2^n}$.

470 Now we discuss how the size of the computable message space $|\text{dom}_{v-u}^\tau|$ can be written
471 when u is one of the children or grandchildren of v .

472 ► **Example 12.** Suppose $u = v_R$. In this case,

$$473 \quad \text{New}_v^\tau(x, y) = \{m|_v : \text{in}^\tau(m_R) = x, y = \text{H}^\tau(m_{RR}) \oplus \text{H}^\tau(m_L) \oplus x_1 \oplus x_2, \\ 474 \quad \quad \quad | (v, (x_1, x_2)) \in \text{dom}(\tau), m(v) = x_1 \oplus \text{H}^\tau(m_L)\}.$$

476 So, $\mathbb{E}_y(|\text{New}_v^\tau(x, y)|) \leq \frac{\text{mc}_{v_R}^\tau(x) \times |\text{ran}_{v_L}^\tau| \times |\tau_v|}{2^n}$, where τ_v denotes the set of elements in the
477 transcript of the form $((v, x), y)$.

478 ► **Example 13.** In the previous case, we could write the expectation of number of newly
479 generated messages in terms of input multi-collision and range size of tree hash. Now, we
480 consider $u = v_{RR}$, i.e. u is a grandchild of v . Refer to Figure 4. Let $h = y \oplus \text{H}^\tau(m_{RRR})$. First,
481 let us look at $|\text{dom}_{v-v_{RR}}^\tau|$.

$$482 \quad |\text{dom}_{v-v_{RR}}^\tau| = \{m|_v : H_1 \oplus h = x'_1 \oplus x'_2, H_2 \oplus y' \oplus h = x''_1 \oplus x''_2, \\ 483 \quad \quad \quad | H_1 = \text{H}^\tau(m_{RL}), H_2 = \text{H}^\tau(m_L), (v_R, (x'_1, x'_2), y'), (v, (x''_1, x''_2), *) \in \tau\}.$$

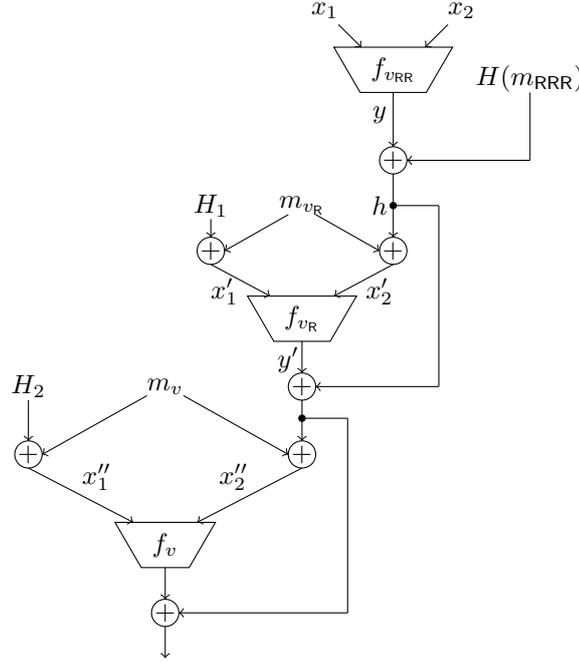
485 Note that this implies $H_1 \oplus H_2 \oplus \bar{y}' = x''_1 \oplus x''_2$, where $\bar{y}' = x'_1 \oplus x'_2 \oplus y'$. Thus,

$$486 \quad |\text{dom}_{v-v_{RR}}^\tau| = \text{mc}(\text{ran}_{v_L}^\tau \oplus \text{ran}_{v_{RL}}^\tau \oplus \bar{f}_{v_R}) \times |\tau_v|,$$

487 where $\bar{f}_{v_R}(u_1, u_2) = u_1 \oplus u_2 \oplus f_{v_R}(u_1, u_2)$. Hence,

$$488 \quad \mathbb{E}_y(|\text{New}_v^\tau(x, y)|) \leq \frac{\text{mc}_{v_{RR}}^\tau(x) \times \text{mc}(\text{ran}_{v_L}^\tau \oplus \text{ran}_{v_{RL}}^\tau \oplus \bar{f}_{v_R}) \times |\tau_v|}{2^n}.$$

489 **ADVERSARY AND ITS QUERIES.** Let \mathcal{L}_v denote the lists of all responses of f_v , for all leaf node
490 v . We can assume that these lists are given to the adversary at the beginning of the game.
491 This is without loss of generality as the inputs to f_v 's have no role in the collision event.
492 However, this is not true for all intermediate nodes (the non-leaf nodes) and so adaptivity
493 of intermediate nodes must be considered. We assume that an adversary makes exactly q



■ **Figure 4** The graph of \mathcal{T}_{v-u} when u is the rightmost grandchild of v .

494 queries to each node. Let $q' := qr$ denote the total number of queries where $r = |V^*|$ and
 495 V^* is the set of non-leaf or intermediate nodes. Let Q_v denote the set of query numbers
 496 for the node v , $v \in V^*$. So for all non-leaf node v , $|Q_v| = q$. Let (x_i, y_i) denote the i -th
 497 query-response pair made to the node v_i . So given transcript τ^{i-1} (transcript after $(i-1)$
 498 queries), the distribution of y is uniform over $\{0, 1\}^n$. For notational simplicity, we use simply
 499 i in a superscript instead of τ^i (the transcript after i -th query) in all above notations defined
 500 so far. For example, $H^i(m)$ denotes the transcript based hash of m where the transcript is τ^i .
 501 We write New_v^i instead of $\text{New}_v^{\tau^{i-1}}(x_i, y_i)$, which represents the set of all newly generated
 502 computable messages at node v immediately after obtaining i -th query-response. We also
 503 ignore the superscript τ^i completely when we all the queries have been made, i.e. $i = q'$. For
 504 example, we write $\text{mc}_v(x)$ instead of $\text{mc}_v^\tau(x)$, when τ is the final transcript, obtained at the
 505 end of all the queries.

- 506 ■ For any computable message m at v , we write $\text{Fin}(m) := i$ to encode the final query
- 507 index after which m is computable.
- 508 ■ For all m for which m_L, m_R are τ -computable, we define $\text{Fin}^*(m) = i$ such that
- 509 $\max\{\text{Fin}(m_L), \text{Fin}(m_R)\} = i$, (i.e. immediately after i -th query the final-input for the
- 510 message m is computable).

511 5.1 Steps of Collision Analysis

512 **PROPER INTERNAL COLLISION.** We say that a **proper internal collision** happens at
 513 $v = (j, b)$ for a transcript τ if for some distinct messages m, m' at v , (i) $H^\tau(m) = H^\tau(m')$, (ii)
 514 $\text{in}^\tau(m) \neq \text{in}^\tau(m')$, and (iii) no collision happens for H_u^τ for all $u \in V(\mathcal{T}_v)$, $u \neq v$. By using
 515 standard reduction, a collision of ABR must have proper internal collision at some node. So
 516 it is sufficient to bound the probability of a proper internal collision at the root node of ABR
 517 as H_v is identical to ABR_s where s denotes the level of the node v . We write $\text{coll} := \text{coll}_l$ to

518 denote the proper internal collision at the root node of \mathcal{T} of height l . The probability of
519 collision of ABR_l can be then bounded as $\sum_{i \leq l} 2^{l-i} \Pr(\text{coll}_i)$.

520 Now, there are two types of collision which can happen for any proper collision at the
521 root. Let us consider the i -th query. This query itself can generate two new computable
522 messages for which the collision occurs. This is the first type of collision. Also, the hash
523 output of one among the new computable messages generated by the i -th query can match
524 with one of the hash outputs generated by the previous queries. We formalize them here:

525 ► **Definition 14** (types of collision). ■ We call a collision pair (M, M') *twin* at the i -th
526 query, $i \in [q']$ if $M, M' \in \text{New}^i$. In this case $\text{in}_{v_i}^i(M) = \text{in}_{v_i}^i(M') = x_i$, where v_i is the
527 node where the i -th query is made.

528 ■ The collision pair is called *non-twin* at the i -th query if exactly one of M and M' is a
529 member of New^i , and the other message is τ^{i-1} -computable.

530 We write coll^i to denote that the proper internal collision happens at the i -th query.
531 Moreover, if it is a twin-collision (or non-twin collision) we denote the event as $\text{coll}^{i,\text{tw}}$ (or
532 $\text{coll}^{i,\text{ntw}}$ respectively). Thus,

$$533 \quad \text{coll} = \bigcup_{i \in [q']} (\text{coll}^{i,\text{ntw}} \cup \text{coll}^{i,\text{tw}}).$$

534 It is easy to see that twin-collision at the root node is not possible as a collision at the right
535 child of the root node is necessary. In notation, $\text{coll}^{i,\text{tw}} = \emptyset$, whenever $v_i = \omega$.

536 5.1.1 Non-Twin Collision Analysis

537 For any non-root, non-leaf node v , we consider cross-collision between H_{-v} and H_ω . Let CC_v^i
538 denote the set of all pairs (m, m') such that (i) m is a complete message, m' is a message
539 for \mathcal{T}_{-v} and (ii) $\text{H}^i(m) = \text{H}_{-v}^i(m')$. Now, a *non-twin collision* can happen at the i -th query
540 (to the node v_i) if freshly generated hash of a message at v_i matches with the v_i -th message
541 block of m' for a cross-collision pair (m, m') of CC_v^{i-1} . Thus,

$$542 \quad \Pr(\text{coll}^{i,\text{ntw}}) \leq \frac{\text{mc}_v^{i-1}(x_i) \times |\text{CC}_v^{i-1}|}{2^n}. \quad (2)$$

543 Now, if $v = \omega$ then the freshly generated hash at the root node is a hash. So, we have,

$$544 \quad \Pr(\text{coll}^{i,\text{ntw}}) \leq \frac{\text{mc}_\omega^{i-1}(x_i) \times L}{2^n}. \quad (3)$$

545 5.1.2 Twin Collision Analysis

546 For any non-root, non-leaf node v and $\delta \in \{0, 1\}^n \setminus \{0^n\}$, let $\text{C}_{\delta,v}$, called δ -collision, denote
547 the set of all pairs (m, m') such that $\text{H}_{-v}^\tau(m) = \text{H}_{-v}^\tau(m')$ and $m(v) \oplus m'(v) = \delta$. We have
548 seen that no twin collision possible at the root node. We define a set

$$549 \quad \Delta^i = \{\text{H}^{i-1}(m_R) \oplus \text{H}^{i-1}(m'_R) : m, m' \in \text{MC}_{v_i}^{i-1}(x_i)\}.$$

550 Now,

$$551 \quad \Pr(\text{coll}^{i,\text{tw}}) \leq \frac{\sum_{\delta \in \Delta} \text{mc}_v^{i-1}(x_i) \times |\text{C}_{\delta,v}^{i-1}|}{2^n}. \quad (4)$$

552 Note that the size of Δ can be at most $(\text{mc}_v^{i-1}(x_i))^2$.

553 Thus, we have seen a collision analysis requires to bound the following random variables.

- 554 1. $\text{mc}_v^{i-1}(x_i)$ for all i (and so for all nodes v),
- 555 2. L : load of the hash,
- 556 3. $|\text{dom}_{-v}^{i-1}|$: load for \mathcal{T}_{-v} which is required to bound the load L ,
- 557 4. $|\text{C}_{\delta,v}^{i-1}|$: size of δ -collision, and
- 558 5. $|\text{CC}_v^{i-1}|$: size of cross-collision.

559 In the following subsection, we present the collision analysis of ABR_2 in which we only need
 560 the input multi-collision and load (which is also bounded in terms of input multi-collision).
 561 We also present a collision analysis of ABR_3 for which the above terms are present.

562 5.2 Collision Analysis of ABR_2 by ABR

563 As discussed above, we can assume that all queries to the compression functions at the leaf
 564 node have been made beforehand and let q denote the number of queries to each oracle. Let
 565 $\mathcal{L}_1, \mathcal{L}_2$ be two lists of outputs of the leaf node functions and let $\omega := (2, 1)$ denote the root
 566 (the only non-leaf node for \mathcal{T} of height 2). Note that the proper collision at height 1 is the
 567 same as the collision of the lists $\mathcal{L}_1, \mathcal{L}_2$. The proper collision at a leaf node can happen with
 568 probability at most $q^2/2^n$.

569 So, we now consider collision at the root $(2, 1)$. For this, we now define a bad event mc_ω
 570 that $\text{mc}_\omega^q > n$. Equivalently, the event can be expressed as $\text{mc}(\mathcal{L}_1 \oplus \mathcal{L}_2) > n$. Note that we
 571 do not have any non-leaf node other than root node. So, the load for hash values L can be
 572 upper bounded as nq , given that mc_ω^q does not hold. Moreover, cross-collision and δ -collision
 573 is also not possible as we do not have any non-leaf, non-root node. Now, it is well known that

$$574 \Pr(\text{mc}(\mathcal{L}_1 \oplus \mathcal{L}_2) > n) \leq \frac{q^2}{2^n}$$

575 (see [1] for details). Thus, the collision probability is bounded by $\frac{(n^2+2)q^2}{2^n}$.

576 5.3 Collision Analysis of $\text{ABR}_h, h \geq 3$ by [1]

577 The proof of [1] is divided into two main parts: (i) bounding the load and (ii) bounding proper
 578 collision probability in terms of the load. ABR fix a parameter ρ (which is chosen to be $n + 1$,
 579 however, the exact value is not relevant to our discussion). Let $L_{i,v} = \sum_{j \leq i, j \in Q_v} |\text{NewH}_v^i|$
 580 represent the total number of generated hash values at v after all i queries. If there is no
 581 collision (which is true while we consider proper internal collision), $L_{i,v}$ is same as the size of
 582 the set $|\text{dom}_v^i|$. To bound load, ABR considered the following bad events (in our notations):

- 583 1. $\text{bad}_{1,v}$: $\text{mc}_v^{q'} > \rho$ at v . Let $\text{bad}_1 := \cup_v \text{bad}_{1,v}$.
- 584 2. $\text{bad}_{2,v}$: $L_{q',v} \geq \rho q$. Let $\text{bad}_2 := \cup_v \text{bad}_{2,v}$.

585 Given $\text{bad}_1, \text{bad}_2$ do not hold, clearly $L \leq 2\rho q$.

586 5.3.1 Step-1: Bounding $\Pr(\text{bad}_1)$

587 Let $\text{bad}_{1,\leq i} = \cup_{(j,b): j \leq i} \text{bad}_{1,(j,b)}$. So it is sufficient to bound $\Pr(\text{bad}_{1,(j,b)} \wedge \neg \text{bad}_{1,< j})$. Let us
 588 fix a query x at $v = (j, b)$. Now, ABR implicitly claimed the following:

589 **Claim 1** [1]: If $\text{MC}_{(j,b)}^{q'}(x) \supseteq \{m_1, \dots, m_\rho\}$ then $\text{in}_{(j-1,2b)}(m_{i,R})$'s are distinct.

590 We note that this claim is not correct. As there can be ρ multi-collision at node $(j-1, 2b)$,
 591 each query can potentially give at most ρ multi-collision at node (j, b) . Hence we can have ρ^2
 592 multi-collision at node (j, b) . Thus, a corrected version of the above claim requires to revise
 593 the parameter ρ depending on the level. So, we may redefine $\text{bad}_{1,(j,b)}$: $\text{mc}_v > \rho^j$ which

594 could solve the issue. This is a fixable minor issue (but will have an impact on the claimed
595 bound).

596 Now to continue with the bound, let us assume that $\text{MC}_v^{q'}(x) \supseteq \{m_1, \dots, m_\rho\}$ such that
597 $\text{in}(m_{i,\text{R}})$'s are distinct and $x = (a, b)$. So we can choose ρ query indices out of q queries to
598 $v_2 := v_{\text{R}}$ in $\binom{q}{\rho}$ ways. For any such choices of ρ tuple $(i_1, i_2, \dots, i_\rho)$ (all queried to v_2), we
599 have

$$600 \quad \Pr(f(x_{i_1}) \oplus \text{H}(m_{1,\text{RR}}) = b, \dots, f(x_{i_\rho}) \oplus \text{H}(m_{\rho,\text{RR}}) = b) = \frac{2\rho q}{2^{n\rho}}$$

601 as there ρq many choices of $\text{H}(m_{i,\text{RR}})$ values (as we assume the load at v_{RR} is less than $2\rho q$).
602 However, the above is true when we consider the cases where $\text{Fin}^*(m_i) = j_i$ where $v_{j_i} = v_2$
603 for all i . The most important case in which the input multi-collision is contributed due to
604 the final queries which are not on right child is not considered in the proof by [1].

605 5.3.2 Step-2: Bounding $\Pr(\text{bad}_2)$

606 Let $\text{bad}_{2,\leq i} = \cup_{(j,b):j\leq i} \text{bad}_{2,v}$. So it is sufficient to bound

$$607 \quad \Pr(\text{bad}_{2,(j,b)} \wedge \neg \text{bad}_{1,<j} \wedge \neg \text{bad}_1 \wedge \neg \text{coll}).$$

608 The main idea to bound the above probability is to bound the expected number of newly
609 generated hash at $v = (j, b)$ over all queries. Then the bad event probability can be bounded
610 by applying Markov's inequality. We have already seen that

$$611 \quad \mathbb{E}_y(|\text{New}_v^i| \mid \tau^{i-1}) = \frac{\text{mc}_{v_i}^\tau(x_i) \times |\text{dom}_{v-v_i}^i|}{2^n}.$$

612 Moreover, we have shown that bounding $|\text{dom}_{v-v_i}^i|$ becomes more complex when v_i is neither
613 v nor a child of v (see Example 13). [1] tried to argue in a different way. ABR showed a
614 bound expectation of load due to all queries of its children (see Example 12). Then, they
615 continued this argument for two levels up (i.e. for the queries on grandchildren as we consider
616 in Example 13). However, they did not analyze this case properly. In particular, they did
617 not consider to bound the $\text{mc}(\text{ran}_{v_{\text{L}}}^\tau \oplus \text{ran}_{v_{\text{RL}}}^\tau \oplus \tilde{f}_{v_{\text{R}}})$. Finally, they claimed the general case
618 by using induction which is clearly unverifiable.

619 5.3.3 Step-3: Proving Collision in terms of Load

620 ABR stated that as analyzed for ABR_2 , given (i) no collision for all primitive, (ii) $\neg \text{bad}_{1,\leq l}$
621 and (iii) $\neg \text{bad}_{2,\leq l}$, the proper internal collision probability at the root node is $\mathbb{E}(L^2)/2^n$
622 where L is the total number computable hash values.

623 There is a fundamental gap in the high level of the proof. As ABR did not explain anything
624 supporting his claim, we show that this statement is not true in general. In particular, we
625 show (in the next subsection) a *hash mode based on 2-to-1 compression function whose load*
626 *is at most q^2 (for any q -query adversary), however, a collision can be found in $O(n)$ queries.*
627 So the above claim cannot be made in general.

628 5.3.4 Missing Step: Twin-Collision Analysis

629 We find that the twin-collision analysis of the ABR hash is missed completely. The bound
630 for δ -collision is not obvious and it requires bounding the probability of some more bad
631 events. In the following section, we have analyzed ABR_3 in which the twin-collision analysis

632 requires a bad event dealing with the multi-collision of xor of random oracle compression
 633 function outputs for two distinct inputs. We do not know any method to bound the number
 634 of cross-collision pairs for a general height tree.

635 5.4 Relationship between Load and Collision Probability

636 A hash function with a high load is unlikely to be collision-resistant. For example,
 637 $\text{xor}(x_1, \dots, x_r) = f_1(x_1) \oplus \dots \oplus f_r(x_r)$ has load 2^r after 2 queries to each oracle f_i . It is easy
 638 to see that the hash function xor is not collision-resistant. Let $r = n$. Then, after making
 639 two queries to each function, we have sufficiently many computable messages. It is then very
 640 easy to find computable collision pairs by solving a linear system of equations. In general, if
 641 the load becomes the order of $2^{n/2}$ then one may expect a collision. However, the converse
 642 need not be true. In other words, we have a hash function where load can not be high, but
 643 still, a collision pair can be generated efficiently.

644 5.4.1 Example of Collision Insecure Hash Functions with Low Load

645 Let MD^f be the MD hash which takes n blocks and initial value is also replaced by
 646 one message block (so exactly $n - 1$ calls of f is required). We define $MD_n^f(M) =$
 647 $MD^{f_1}(M) \parallel \dots \parallel MD^{f_n}(M)$ which is n^2 -to- n^2 hash function. Now we define a hash function
 648 $H(M_1, M_2)$ for $M_1, M_2 \in \{0, 1\}^{n^2}$:

- 649 1. Let $(C_1, C_2) = (MD_n^f(M_1) \oplus M_2, MD_n^g(C_1) \oplus M_1)$ (two round LR construction which is
 650 invertible).
- 651 2. Let h_1, \dots, h_n be $2n$ -to- n functions. The final hash output is defined as $h_1(x_1) \oplus \dots \oplus$
 652 $h_n(x_n)$ where $C_1 \parallel C_2 = x_1 \parallel \dots \parallel x_n$, $x_i \in \{0, 1\}^{2n}$.

653 Note that we cannot compute (C_1, C_2) for more than q^2 messages assuming there is no
 654 collisions in f and g functions. So, $L(q) \leq q^2$ for any q -query adversary.

655 A COLLISION ATTACK. Now, we construct a collision finding algorithm for the above hash.
 656 It first finds collision pair for xor function $h_1 \oplus \dots \oplus h_n$ (can be achieved easily by making $2n$
 657 queries altogether). Let (C, C') be a collision pair. We can easily invert C and C' to obtain
 658 M and M' respectively. Clearly, (M, M') is a computable collision pair.

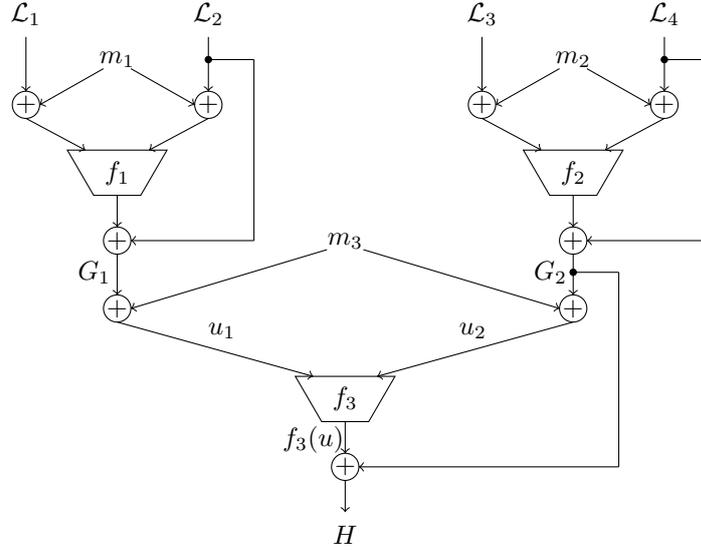
659 6 Analysis of ABR of height 3

660 In this section, we show that the ABR_3 construction achieves birthday security. In particular,
 661 we prove the following theorem:

662 ► **Theorem 15** (collision theorem for ABR_3). *For any adversary \mathcal{A} making at most q queries
 663 to each compression function modeled to be random oracle, we have*

$$664 \quad \text{Adv}_{H^f}^{\text{coll}}(\mathcal{A}) \leq \frac{6n^5q^2 + 3n^4q^2 + 2n^4q + 2n^2q^2 + 13q^2}{2^n}. \quad (5)$$

666 Let $\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3, \mathcal{L}_4$ be the four lists of size q each corresponding to the outputs of
 667 $f_{1,1}, f_{1,2}, f_{1,3}, f_{1,4}$ respectively. We can assume that these lists are given to the adver-
 668 sary at the beginning of the game. This is without loss of generality as the inputs to $f_{1,i}$'s
 669 are independent from the rest of the transcripts. Also, for ease of notation, from now on
 670 we denote $f_{2,1}$ by f_1 , $f_{2,2}$ by f_2 and $f_{3,1}$ by f_3 . If the input to any of the functions is
 671 $u = (u_1, u_2)$, we define $u^\oplus = u_1 \oplus u_2$. Also, if $f_3(u) = v$, then we define $\bar{f}_3(u) = u^\oplus \oplus v$. As
 672 f_3 is a random oracle, the output distributions of \bar{f}_3 are uniform and independent. Let Q_j



■ **Figure 5** ABR_3 according to our new notation when the query $u = (u_1||u_2)$ is made to f_3 .

673 be the set of queries to f_j . We assume $|Q_j| = q$ for $j = 1, 2, 3$. Also, let Q_j^i denote the set of
 674 queries to Q_j up to the i -th query (including the i -th one). Let $\mathcal{G}_1 = O_{(2,1)}$ denote the set of
 675 intermediate hash outputs at node $(2, 1)$ and $\mathcal{G}_2 = O_{(2,2)}$. Let \mathcal{H} denote the set of final hash
 676 outputs of ABR_3 . Refer to Figure 5 for a pictorial representation. We follow the general
 677 approach as described before. We have already shown the collision bound for ABR_2 and so it
 678 is sufficient to bound proper collision at the root for ABR_3 .

679 As we have seen above, the collision analysis requires us to bound some random variables.
 680 We first define some bad events to bound these random variables.

681 ► **Definition 16** (list collision). *The first bad event we consider is:*

682 ■ B_0 : There exists a collision in at least one of the lists $\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3, \mathcal{L}_4, \{f_1(u) : u \in Q_1\},$
 683 $\{f_2(u) : u \in Q_2\}, \{f_3(u) : u \in Q_3\}$.

684 Since f is modeled as a random function, the collision probability in any of the lists is at
 685 most $q^2/2^n$. Hence, $\Pr(B_0) \leq 7q^2/2^n$.

686 ► **Definition 17** (bad event on input multi-collision). *We define the following bad events:*

687 ■ B_1 : $\text{mc}(\mathcal{L}_1 \oplus \mathcal{L}_2) > n$, or $\text{mc}(\mathcal{L}_3 \oplus \mathcal{L}_4) > n$,

688 ■ B_2 : $\text{mc}(\mathcal{G}_1 \oplus \mathcal{G}_2) > n^2$.

689 We now state some simple observations related to input multi-collision:

- 690 1. Given that B_1 does not hold, $\text{mc}_{(2,1)}, \text{mc}_{(2,2)} \leq n$ and so $|\mathcal{G}_1|, |\mathcal{G}_2| \leq nq$.
- 691 2. Given that B_2 does not hold, $\text{mc}_{(3,1)} \leq n^2$ and so $L_{(3,1)} \leq n^2q$.
- 692 3. Note, $|\text{dom}(\mathcal{T}_{-(2,1)})|, |\text{dom}(\mathcal{T}_{-(2,2)})| \leq nq^2$. So, $\mathbb{E}(L_{(2,1)}), \mathbb{E}(L_{(2,2)}) \leq n^2q^3/2^n$. By
 693 Markov's inequality, $\Pr(L > 3n^2q) \leq 2q^2/2^n$ ($3n^2q$ because we include $L_{(3,1)}$ as well).
- 694 4. By using a similar argument as we applied for multi-collision, we have $\Pr(B_1) \leq 2q^2/2^n$.
- 695 5. Now, given that B_1 does not hold and B_2 holds, there must exist at least n distinct inputs
 696 to f_2 leading to n^2 input multi-collision. So, we can similarly prove $\Pr(B_2) \leq q^2/2^n$.

697 We say that bad_{mc} holds if either B_1 or B_2 happens, or $L > 3n^2q$. Then, from above,
 698 $\Pr(\text{bad}_{mc}) \leq 3q^2/2^n$. We now define bad events which would be used to bound cross-collision.

699 ► **Definition 18** (bad event on cross-collision). *We define the following bad events:*

11:20 Revisiting Collision and Local Opening Analysis of ABR Hash

700 ■ $B3: |\{(G_2, f_3(u), H) : G_2 \oplus f_3(u) \oplus H = 0; G_2 \in \mathcal{G}_2, u \in Q_3, H \in \mathcal{H}\}| > 3n^4q.$

701 ■ $B4: |\{(G_1, f_3(u), H) : G_1 \oplus f_3(u) \oplus H = 0; G_1 \in \mathcal{G}_1, u \in Q_3, H \in \mathcal{H}\}| > 3n^4q.$

702 *We say that bad_{cc} holds if any one of the above happens.*

703 If the i -th query is made at f_2 , an intermediate hash output G_2 generated at this level
 704 due to this query can match with a query u already done to f_3 to generate a final hash
 705 output H which was already previously generated by the first $i - 1$ queries. The event $B3$
 706 implies that the number of such triplets $(G_2, f_3(u), H)$ is more than $3n^4q$. $B4$ has a similar
 707 implication when we consider \mathcal{G}_1 instead of \mathcal{G}_2 .

708 ► **Lemma 19.** $\Pr(\text{bad}_{cc} \wedge \neg \text{bad}_{mc}) \leq 2q^2/2^n.$

709 **Proof.** $\Pr(\text{mc}(\mathcal{G}_2 \oplus \text{ran}(f_3)) > n^2) \leq q^2/2^n.$ The proof is similar to that of event $B2$. Hence,
 710 for a fixed $H \in \mathcal{H}$, we have

711 $\Pr[|\{(G_2, f_3(u), H) : G_2 \oplus f_3(u) \oplus H = 0; G_2 \in \mathcal{G}_2, u \in Q_3\}| > n^2] \leq q^2/2^n.$

712 Now, there are $3n^2q$ choices for H . Therefore, $\Pr(B3 \wedge \neg \text{bad}_{mc}) \leq q^2/2^n.$ A similar argument
 713 works for $B4$. Hence,

714 $\Pr(\text{bad}_{cc} \wedge \neg \text{bad}_{mc}) \leq 2q^2/2^n. \quad \blacktriangleleft$

715 Given that bad_{cc} does not hold, $|\text{CC}_{(2,1)}| \leq 3n^4q$ (or $|\text{CC}_{(2,2)}| \leq 3n^4q$ respectively). We
 716 finally define bad events which would be used to bound δ -collision pairs.

717 ► **Definition 20** (bad event on δ -collision). *We define the following bad event:*

718 ■ $B5: \text{mc}(\bar{f}_3(u) \oplus \bar{f}_3(u')) > n.$

719 *We say that bad_δ holds if the above happens.*

720 ► **Lemma 21.** $\Pr(\text{bad}_\delta) \leq \frac{q^2}{2^n}.$

721 **Proof.** Since $f_3(u)$ is random, $\bar{f}_3(u) = f_3(u) \oplus u^\oplus$ is also random. Therefore, bounding $B5$
 722 is similar to bounding $B1$. \blacktriangleleft

723 Given that bad_δ does not hold, $|\text{C}_\delta| \leq n$. Let $\text{bad} = B0 \cup \text{bad}_{mc} \cup \text{bad}_{cc} \cup \text{bad}_\delta$. Then,
 724 $\Pr(\text{bad}) \leq \frac{13q^2}{2^n}.$

725 6.1 Collision Analysis

726 We assume that bad does not hold. Since $\text{coll} = \bigcup_{i \in [q], v \in V \setminus \mathcal{L}} (\text{coll}_v^{i,\text{ntw}} \cup \text{coll}_v^{i,\text{tw}})$, we need to

727 bound $\text{coll}_v^{i,\text{ntw}}$ and $\text{coll}_v^{i,\text{tw}}$ for $v = (2,1), (2,2), (3,1)$. In the following lemmas, we bound
 728 them, assuming bad does not occur. We already know that $\text{coll}_{(3,1)}^{i,\text{tw}}$ does not occur.

729 ► **Lemma 22.** $\Pr(\text{coll}_{(3,1)}^{i,\text{ntw}} | \neg \text{bad}) \leq \frac{3n^4q}{2^n}.$

730 **Proof.** As seen above in equation 3, $\Pr(\text{coll}_{(3,1)}^{i,\text{ntw}}) \leq \frac{\text{mc}_{(3,1)}^{i-1}(x_i) \times L}{2^n}.$

731 Given $\neg \text{bad}$, $\text{mc}_{(3,1)} \leq n^2$ and $L \leq 3n^2q$. Hence, $\Pr(\text{coll}_{(3,1)}^{i,\text{ntw}} | \neg \text{bad}) \leq \frac{3n^4q}{2^n}. \quad \blacktriangleleft$

732 ► **Lemma 23.** $\Pr(\text{coll}_{(2,1)}^{i,\text{ntw}} | \neg \text{bad}) \leq \frac{3n^5q}{2^n}.$

733 **Proof.** As seen above in equation 3, $\Pr(\text{coll}_{(2,1)}^{i,\text{ntw}}) \leq \frac{\text{mc}_{(2,1)}^{i-1}(x_i) \times |\text{CC}_{(2,1)}^{i-1}|}{2^n}$.

734 Given $\neg\text{bad}$, $\text{mc}_{(2,1)} \leq n$ and $|\text{CC}_{(2,1)}| \leq 3n^4q$. Hence, $\Pr(\text{coll}_{(2,1)}^{i,\text{ntw}} | \neg\text{bad}) \leq \frac{3n^5q}{2^n}$. ◀

735 ▶ **Lemma 24.** $\Pr(\text{coll}_{(2,2)}^{i,\text{ntw}} | \neg\text{bad}) \leq \frac{3n^5q}{2^n}$.

736 **Proof.** This proof is similar to that of the previous lemma.

737 $\Pr(\text{coll}_{(2,2)}^{i,\text{ntw}}) \leq \frac{\text{mc}_{(2,2)}^{i-1}(x_i) \times |\text{CC}_{(2,2)}^{i-1}|}{2^n}$.

738 Given $\neg\text{bad}$, $\text{mc}_{(2,2)} \leq n$ and $|\text{CC}_{(2,2)}| \leq 3n^4q$. Hence, $\Pr(\text{coll}_{(2,2)}^{i,\text{ntw}} | \neg\text{bad}) \leq \frac{3n^5q}{2^n}$. ◀

739 ▶ **Lemma 25.** $\Pr(\text{coll}_{(2,1)}^{i,\text{tw}} | \neg\text{bad}) \leq \frac{n^4}{2^n}$.

740 **Proof.** As seen above in equation 4, $\Pr(\text{coll}_{(2,1)}^{i,\text{tw}}) \leq \frac{\sum_{\delta \in \Delta} \text{mc}_{(2,1)}^{i-1}(x_i) \times |\mathbf{C}_{\delta,(2,1)}^{i-1}|}{2^n}$.

741 Given $\neg\text{bad}$, $\text{mc}_{(2,1)} \leq n$, $|\Delta| \leq (\text{mc}_{(2,1)}^{i-1}(x_i))^2 \leq n^2$ and $|\mathbf{C}_{\delta,(2,1)}| \leq n$. Hence,

742 $\Pr(\text{coll}_{(2,1)}^{i,\text{tw}} | \neg\text{bad}) \leq \frac{n^4}{2^n}$. ◀

743 ▶ **Lemma 26.** $\Pr(\text{coll}_{(2,2)}^{i,\text{tw}} | \neg\text{bad}) \leq \frac{n^4}{2^n}$.

744 **Proof.** This proof is similar to that of the previous lemma.

745 $\Pr(\text{coll}_{(2,2)}^{i,\text{tw}}) \leq \frac{\sum_{\delta \in \Delta} \text{mc}_{(2,2)}^{i-1}(x_i) \times |\mathbf{C}_{\delta,(2,2)}^{i-1}|}{2^n}$.

746 Given $\neg\text{bad}$, $\text{mc}_{(2,2)} \leq n$, $|\Delta| \leq (\text{mc}_{(2,2)}^{i-1}(x_i))^2 \leq n^2$ and $|\mathbf{C}_{\delta,(2,2)}| \leq n$. Hence,

747 $\Pr(\text{coll}_{(2,2)}^{i,\text{tw}} | \neg\text{bad}) \leq \frac{n^4}{2^n}$. ◀

748 From the above lemmas, we have

749
$$\Pr(\text{coll} | \neg\text{bad}) \leq \sum_{i \in [q], v \in V \setminus \mathcal{L}} \Pr(\text{coll}_v^{i,\text{ntw}} | \neg\text{bad}) + \Pr(\text{coll}_v^{i,\text{tw}} | \neg\text{bad}) \leq \frac{6n^5q^2 + 3n^4q^2 + 2n^4q}{2^n}.$$

750 Therefore, $\Pr(\text{coll}) \leq \Pr(\text{coll} | \neg\text{bad}) + \Pr(\text{bad}) \leq \frac{6n^5q^2 + 3n^4q^2 + 2n^4q + 13q^2}{2^n}$.

751 Note that we have bound the proper collision probability at the root for ABR_3 . Since
752 B_0 does not occur, collision does not occur at the leaf node. As seen in section 5.2, the
753 probability that proper collision occurs at node $(2, 1)$ (resp. $(2, 2)$) is bounded above by $\frac{n^2q^2}{2^n}$.

754 Hence, the theorem is proved.

755 **7 Conclusion**

756 In this paper, we revisit the collision security of the ABR hash. We found that there is a
757 serious gap in the analysis of collision security. Some missing and important cases have also
758 been identified. In this paper, we have shown collision security for level 3. Several new bad
759 events have been identified in ABR_3 which were not considered for the general hash. We
760 leave the collision security analysis open for general hash. Thus, the optimality of Stam's
761 bound remains open for an arbitrary domain hash.

762 We have also found that the ABR hash cannot have any non-trivial local opening which
763 can give birthday bound security. This shows a limitation in terms of applications in local
764 opening. In particular, the efficient local opening proposed by [1] can be broken in $O(2^{n/2l})$
765 query complexity.

766 — **References** —

- 767 **1** Elena Andreeva, Rishiraj Bhattacharyya, and Arnab Roy. Compactness of hashing modes and
768 efficiency beyond merkle tree. In *Advances in Cryptology - EUROCRYPT 2021 - 40th Annual*
769 *International Conference on the Theory and Applications of Cryptographic Techniques, Zagreb,*
770 *Croatia, October 17-21, 2021, Proceedings, Part II*, pages 92–123. Springer, 2021.
- 771 **2** John Black, Martin Cochran, and Thomas Shrimpton. On the impossibility of highly-efficient
772 blockcipher-based hash functions. In *EUROCRYPT*, volume 3494 of *Lecture Notes in Computer*
773 *Science*, pages 526–541. Springer, 2005.
- 774 **3** Ivan Damgård. A design principle for hash functions. In *CRYPTO*, volume 435 of *Lecture*
775 *Notes in Computer Science*, pages 416–427. Springer, 1989.
- 776 **4** Yevgeniy Dodis, Dmitry Khovratovich, Nicky Mouha, and Mridul Nandi. T5: Hashing five
777 inputs with three compression calls. In *2nd Conference on Information-Theoretic Cryptography,*
778 *ITC 2021, July 23-26, 2021, Virtual Conference*, pages 24:1–24:23. Schloss Dagstuhl - Leibniz-
779 Zentrum für Informatik, 2021.
- 780 **5** Bart Mennink and Bart Preneel. Efficient parallelizable hashing using small non-compressing
781 primitives. *Int. J. Inf. Sec.*, 15(3):285–300, 2016.
- 782 **6** Ralph C. Merkle. Protocols for public key cryptosystems. In *IEEE Symposium on Security*
783 *and Privacy*, pages 122–134. IEEE Computer Society, 1980.
- 784 **7** Ralph C. Merkle. One way hash functions and DES. In *CRYPTO*, volume 435 of *Lecture*
785 *Notes in Computer Science*, pages 428–446. Springer, 1989.
- 786 **8** Phillip Rogaway and John P. Steinberger. Constructing cryptographic hash functions from
787 fixed-key blockciphers. In *CRYPTO*, volume 5157 of *Lecture Notes in Computer Science*, pages
788 433–450. Springer, 2008.
- 789 **9** Phillip Rogaway and John P. Steinberger. Security/efficiency tradeoffs for permutation-based
790 hashing. In *EUROCRYPT*, volume 4965 of *Lecture Notes in Computer Science*, pages 220–236.
791 Springer, 2008.
- 792 **10** Thomas Shrimpton and Martijn Stam. Building a collision-resistant compression function
793 from non-compressing primitives. In *ICALP (2)*, volume 5126 of *Lecture Notes in Computer*
794 *Science*, pages 643–654. Springer, 2008.
- 795 **11** Martijn Stam. Beyond uniformity: Better security/efficiency tradeoffs for compression functions.
796 In *Advances in Cryptology - CRYPTO 2008, 28th Annual International Cryptology Conference,*
797 *Santa Barbara, CA, USA, August 17-21, 2008. Proceedings*, pages 397–412. Springer, 2008.
- 798 **12** John P. Steinberger. Stam’s collision resistance conjecture. In *EUROCRYPT*, volume 6110 of
799 *Lecture Notes in Computer Science*, pages 597–615. Springer, 2010.
- 800 **13** John P. Steinberger, Xiaoming Sun, and Zhe Yang. Stam’s conjecture and threshold phenomena
801 in collision resistance. In *CRYPTO*, volume 7417 of *Lecture Notes in Computer Science*, pages
802 384–405. Springer, 2012.
- 803 **14** David A. Wagner. A generalized birthday problem. In *CRYPTO*, volume 2442 of *Lecture*
804 *Notes in Computer Science*, pages 288–303. Springer, 2002.