

# Verification of the CVM Algorithm with a Functional Probabilistic Invariant

Emin Karayel ✉ 

School of Computation, Information and Technology, Technical University of Munich, Germany

Seng Joe Watt ✉ 

Institute for Infocomm Research (I<sup>2</sup>R), A\*STAR, Singapore

Derek Khu ✉ 

Institute for Infocomm Research (I<sup>2</sup>R), A\*STAR, Singapore

Kuldeep S. Meel ✉ 

Georgia Institute of Technology, Atlanta, GA, USA

University of Toronto, Canada

Yong Kiam Tan ✉ 

Institute for Infocomm Research (I<sup>2</sup>R), A\*STAR, Singapore

Nanyang Technological University, Singapore

---

## Abstract

Estimating the number of distinct elements in a data stream is a classic problem with numerous applications in computer science. We formalize a recent, remarkably simple, randomized algorithm for this problem due to Chakraborty, Vinodchandran, and Meel (called the CVM algorithm). Their algorithm deviated considerably from the state of the art, due to its avoidance of intricate derandomization techniques, while still maintaining a close-to-optimal logarithmic space complexity.

Central to our formalization is a new proof technique based on functional probabilistic invariants, which allows us to derive concentration bounds using the Cramér–Chernoff method without relying on independence. This simplifies the formal analysis considerably compared to the original proof by Chakraborty et al. Moreover, our technique opens up the possible algorithm design space; we demonstrate this by introducing and verifying a new variant of the CVM algorithm that is both total and unbiased – neither of which is a property of the original algorithm. In this paper, we introduce the proof technique, describe its use in mechanizing both versions of the CVM algorithm in Isabelle/HOL, and present a supporting formalized library on negatively associated random variables used to verify the latter variant.

**2012 ACM Subject Classification** Theory of computation → Logic and verification; Theory of computation → Higher order logic; Mathematics of computing → Probabilistic algorithms; Theory of computation → Pseudorandomness and derandomization

**Keywords and phrases** Verification, Isabelle/HOL, Randomized Algorithms, Distinct Elements

**Digital Object Identifier** 10.4230/LIPIcs.ITP.2025.34

## Supplementary Material

*Software:* [https://isa-afp.org/entries/CVM\\_Distinct\\_Elements.html](https://isa-afp.org/entries/CVM_Distinct_Elements.html) [34]

*Software:* [https://github.com/joewatt95/CVM/tree/main/isabelle/CVM\\_Transforms](https://github.com/joewatt95/CVM/tree/main/isabelle/CVM_Transforms) [48]

archived at [swh:1:dir:7df92f0347e7bd150efef42e3edbbde0037498cc](https://swh.1:dir:7df92f0347e7bd150efef42e3edbbde0037498cc)

*Software:* [https://isa-afp.org/entries/Negative\\_Association.html](https://isa-afp.org/entries/Negative_Association.html) [33]

**Funding** *Seng Joe Watt:* Singapore NRF Fellowship Programme NRF-NRFF16-2024-0002.

*Kuldeep S. Meel:* Natural Sciences and Engineering Research Council of Canada (NSERC), funding reference [RGPIN-2024-05956].

*Yong Kiam Tan:* Singapore NRF Fellowship Programme NRF-NRFF16-2024-0002.



© Emin Karayel, Seng Joe Watt, Derek Khu, Kuldeep S. Meel, and Yong Kiam Tan; licensed under Creative Commons License CC-BY 4.0

16th International Conference on Interactive Theorem Proving (ITP 2025).

Editors: Yannick Forster and Chantal Keller; Article No. 34; pp. 34:1–34:20

Leibniz International Proceedings in Informatics



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

## 1 Introduction

In 2022, Chakraborty, Vinodchandran, and Meel [10] published a remarkable streaming algorithm for the distinct elements problem [41]. Indeed, Knuth later wrote a note on the algorithm [35], pointing out its interesting properties and christening it the *CVM* algorithm (which we use for the rest of this paper). One striking property of the CVM algorithm is that, in contrast to every other known algorithm for the problem, it does not rely on hashing the stream elements. Instead, the algorithm could theoretically be implemented in a setting where objects in the data stream only allow for equality comparisons. Another property is its simplicity, which is why the authors called it “an algorithm for the textbook”. The algorithm is displayed in its entirety in Algorithm 1.

■ **Algorithm 1** CVM algorithm for distinct elements estimation [10].

**Input:** Stream elements  $a_1, \dots, a_l$ ,  $0 < \varepsilon$ ,  $0 < \delta < 1$ .

**Output:** A cardinality estimate  $R$  for set  $A = \{a_1, \dots, a_l\}$  such that  $\mathcal{P}(|R - |A|| > \varepsilon|A|) \leq \delta$

```

1:  $\chi \leftarrow \{\}, p \leftarrow 1, n = \lceil \frac{12}{\varepsilon^2} \ln(\frac{6l}{\delta}) \rceil$ 
2: for  $i \leftarrow 1$  to  $l$  do
3:    $b \xleftarrow{\$} \text{Ber}(p)$  ▷ randomly sample a bit  $b$  from the Bernoulli distribution
4:   if  $b$  then ▷ insert  $a_i$  if  $b$  is true (with prob.  $p$ )
5:      $\chi \leftarrow \chi \cup \{a_i\}$ 
6:   else ▷ remove  $a_i$  otherwise
7:      $\chi \leftarrow \chi - \{a_i\}$ 
8:   if  $|\chi| = n$  then ▷ if buffer  $\chi$  is full
9:      $\chi \xleftarrow{\$} \text{subsample}(\chi)$  ▷ discard elements of  $\chi$  independently with prob.  $\frac{1}{2}$ 
10:     $p \leftarrow \frac{p}{2}$ 
11:   if  $|\chi| = n$  then return  $\perp$  ▷ fail if  $\chi$  remains full
12: return  $\frac{|\chi|}{p}$  ▷ estimate cardinality of  $A$ 

```

The pen-and-paper analysis of CVM [10, 11] relies on a sequence of transformations of the algorithm. The reason for these transformations is that standard methods for analyzing randomized algorithms, such as Chernoff–Hoeffding bounds, usually make statements about independent random variables. However, for Algorithm 1, the state variables are far from being independent.<sup>1</sup> For example, in Line 3 the Bernoulli distribution is sampled with the parameter  $p$ , which itself depends on previous random operations; similarly, the subsampling step in Line 9 is only applied if the buffer  $\chi$  is full, which also depends on previous random operations. The aforementioned sequence of transformations by Chakraborty et al. results in another randomized algorithm which can be analyzed using standard methods, and from which the desired results for the original algorithm can be deduced. To our knowledge, it seems impossible to analyze Algorithm 1 more directly using known textbook methods [3, 39, 40].

In this paper, we present a new technique for analyzing randomized algorithms which yields a direct and substantially more general proof of the CVM algorithm. Our approach is very similar to how deterministic algorithms are verified using loop invariants. The key difference is that our choice of “loop invariant” for the randomized streaming algorithm is a functional probabilistic inequality, namely, we consider invariants of the form:

$$\mathbb{E}[h] \leq h(c)$$

<sup>1</sup> There is an incorrect claim in the initial published proof of CVM [10, Claim 6] that the indicator functions for elements in  $\chi$  are independent; a later version by the same authors [11] provides a correct proof. The original error serves as a side motivation for this work.

where the expectation is taken over the distribution of the state of the algorithm;  $h$  is allowed to range over a class of functions (mapping states to real values); and  $c$  is a fixed state (possibly chosen differently at each loop iteration). By first establishing such an invariant for Algorithm 1, we can then use it (via different choices of  $h$ ) to establish error bounds for the algorithm. We coined the term “functional probabilistic invariant”, borrowing loosely from the theory of the calculus of variations, where scalar-valued maps – called functionals – from the problem space are used to solve optimization problems. In our case we are using scalar-valued maps – i.e., functionals – from the distribution of the state of the algorithm, which lead us to the name for our new technique. We believe the new proof remains accessible at the undergraduate level, albeit with some exposure to mechanized theorem proving.

To show the generality of our technique, we introduce a new variant of the CVM algorithm, where the subsampling step in Line 9 of Algorithm 1 selects a random  $m$ -subset of  $\chi$  instead of independently discarding each element with some probability. This variant has the benefit that it is *total* (never returns  $\perp$ ) because the second check in Line 11 becomes obsolete. More interestingly, the variant is *unbiased*, i.e., the expected value of the algorithm’s output is exactly the cardinality of the elements in the stream; this is a new property that neither the original CVM algorithm nor classic algorithms for the distinct elements problem possess.

The modified subsampling step leads to additional dependence for the elements in  $\chi$  which cannot be readily removed using transformations as was done in the original proof. Instead, we verify the new variant with our probabilistic invariant-based approach, using results from the theory of negatively associated random variables [29] to establish the desired functional invariant. The concept of negative association is a generalization of independence; importantly, negatively associated variables observe closure properties and fulfill Chernoff–Hoeffding bounds similarly to independent random variables. It should be stressed that the theory of negative association is orthogonal to our new technique, but its formalization is also a contribution of this work.

In summary, our main contributions are:

- Introduction of a new technique using functional probabilistic invariants to verify tail-bounds for randomized algorithms inductively/recursively.
- Verification of the original CVM algorithm using our new technique.
- Presentation and verification of a new variant of CVM that is total and unbiased.
- Formalization of a theory of negatively associated random variables used to analyze the new CVM variant.

We carried out the mechanizations using Isabelle/HOL [42], which comes with a large repository of foundational libraries [1] for the verification of randomized algorithms. We have also mechanized the transformation-based CVM proof by Chakraborty et al. [10, 11], which provides a rough point of comparison: verification of the CVM algorithm using our new technique required only 1003 lines, while the original proof required 2634 lines.<sup>2</sup>

The rest of this paper is organized as follows. Section 2 provides background information on randomized algorithms, in particular on their semantics in Isabelle/HOL. Section 3 introduces our new technique and explains how probabilistic loop invariants can be used to establish tail bounds for the original CVM algorithm. Section 4 introduces the concept of negative association and our new total and unbiased variant of the CVM algorithm. Section 5 presents the formalization of both variants of the algorithm, and Section 6 describes our

<sup>2</sup> We count the total number of lines of Isabelle code in the whole project, excluding empty, comment, and presentation-related lines.

new formalized library on negatively associated random variables. Section 7 discusses some challenges faced in our alternative verification of CVM using the transformation-based proof by Chakraborty et al. The final sections present related work and a summary of our results.

The supplementary material contains:

- formalization [34] of the CVM algorithm, both the original version (Algorithm 1) and our new version (Algorithm 3) using functional probabilistic invariants;
- formalization [33] of a library for negative association; and
- formalization [48] of the CVM algorithm following the proof by Chakraborty et al. [10, 11].

## 2 Background

### 2.1 Randomized Algorithms and Distinct Elements

The CVM algorithm is a *streaming* algorithm for the distinct elements problem. As shown in Algorithm 1, given a data stream  $a_1, \dots, a_l$ , the goal of such algorithms is to return an accurate cardinality estimate for the set  $A = \{a_1, \dots, a_l\}$ .

Importantly, CVM is a *probably approximately correct* (PAC) algorithm where its output estimate  $R$  satisfies  $\mathcal{P}(|R - |A|| > \varepsilon|A|) \leq \delta$  for parameters  $\varepsilon$  and  $\delta$ , i.e., the probability that the relative error of  $R$  with respect to  $|A|$  exceeds  $\varepsilon$  is at most  $\delta$ . Moreover, let us assume that the space needed to store each element in the stream is  $b$  bits, then the CVM algorithm requires only  $\mathcal{O}(\varepsilon^{-2} b \ln(\delta^{-1} l))$  bits of mutable state, which is far less than storing each stream element deterministically.

► **Remark 1.** The asymptotically optimal randomized algorithm for distinct elements requires  $\mathcal{O}(\varepsilon^{-2} \ln \delta + b)$  bits, but it requires more advanced algorithmic techniques. It would not be possible to present using such elementary steps as in Algorithm 1 as it involves computations in finite fields and random walks in expander graphs [8, 32].  $\lrcorner$

### 2.2 Semantics of Randomized Algorithms

We briefly review how reasoning about randomized algorithms works in Isabelle/HOL using the Giry monad [21]. Multiple authors provide more thorough discussions of the concept in the context of Isabelle and other proof assistants [4, 17, 37].

The key idea is to model a randomized algorithm as a probability space representing the distribution of its results. As an example, let us consider Algorithm 2.

■ **Algorithm 2** Example for sequential composition.

---

```

1:  $p \stackrel{\$}{\leftarrow} \text{Ber}(\frac{1}{2})$ 
2:  $q \stackrel{\$}{\leftarrow} \text{Ber}(\frac{1}{3} + \frac{p}{2})$ 
3: return  $q$ 

```

---

In the first step, Algorithm 2 flips a fair coin, such that  $p$  is 1 with probability  $\frac{1}{2}$  and 0 otherwise; the notation  $\text{Ber}(p)$  represents the Bernoulli distribution. In the second step, the algorithm flips a coin  $q$  which depends on  $p$ . This has the consequence that, to semantically model  $q$ , we have to consider functions returning probability spaces, like:  $p \mapsto \text{Ber}(\frac{1}{3} + \frac{p}{2})$ , which is being *bound* to the distribution of  $p$ . The resulting distribution for  $q$  is a *compound distribution* resulting from a combination of  $\text{Ber}(\frac{1}{3})$  (when  $p = 0$ ) and  $\text{Ber}(\frac{5}{6})$  (when  $p = 1$ ).

This example captures the main aspects of modeling randomized algorithms in the Giry monad. Indeed, randomized algorithms can be modeled using the following ingredients:

**Primitive Random Operations.** For example, a simple fair coin flip is represented using the Bernoulli distribution,  $\text{Ber}(\frac{1}{2})$ .

**Return Combinator.** Given an element  $x$ , we can construct the singleton probability space, assigning probability 1 to  $x$  and 0 to everything else. In monad notation, this is written as: **return**  $x$ .

**Bind Combinator.** The bind combinator represents sequential composition of two randomized algorithms  $m$  and  $f$ , where the latter randomized algorithm consumes the output of the former; in monad notation, this is:  $m \gg f$ . Mathematically, this is the most involved operation, because  $f$  is a function returning probability spaces, which takes inputs from the probability space  $m$ .

Let us consider an event  $A$  in the probability space  $m \gg f$ . Its probability can be evaluated by integrating over its probabilities in  $f$  with respect to  $m$ :

$$\mathcal{P}_{m \gg f}(A) = \int_m \mathcal{P}_{f(x)}(A) dx.$$

Another key property is the calculation of expectations; if  $h$  is a random variable over  $m \gg f$ , we can compute its expectation as:

$$\mathbb{E}_{m \gg f}[h] = \int_m \mathbb{E}_{f(x)}[h] dx. \quad (1)$$

Equation 1 is crucially used to establish the invariants we introduce next in Section 3.

### 3 Functional Probabilistic Invariants

In this section, we will derive our new technique using Algorithm 1 as an example. Let us start by briefly reviewing the algorithm – its state is a buffer  $\chi$  (initially empty) and a fraction  $p > 0$  (initially set to 1). The buffer tracks a subset of the elements of the stream encountered so far, with maximal size  $n$  chosen according to the desired accuracy parameters  $\varepsilon$ ,  $\delta$ , and the stream size  $l$ . The algorithm iterates over the stream elements, adding each one to the buffer with probability  $p$  or conversely – if the current stream element is already in the buffer – removing it with probability  $(1 - p)$  (Lines 3–7). If the number of elements in the buffer reaches the maximal size  $n$ , the subsampling operation is executed, which discards each element in  $\chi$  independently with probability  $\frac{1}{2}$ ; then,  $p$  is adjusted to reflect the fact that the buffer now contains each element with probability  $p_{\text{new}} = \frac{p_{\text{old}}}{2}$  (Lines 8–10). If the subsampling operation fails, i.e., if no elements get discarded, then the algorithm fails returning  $\perp$  (Line 11). After processing the stream, the algorithm returns  $\frac{|\chi|}{p}$  as a probably-approximately correct estimate for the number of distinct elements in the stream.

► **Remark 2.** For our discussion below, it is convenient to analyze Algorithm 1 without Line 11, i.e., we will skip the second check of  $|\chi| = n$  determining whether the subsampling step succeeded. This modified version simplifies our analysis as we do not have to worry about the possibility of the algorithm failing (returning  $\perp$ ). This transformation is also used in the original CVM proof [11], where the total variational distance between these two variants of the algorithms is shown to be at most  $\frac{\delta}{2}$ . Thus, probability bounds derived for the modified version can be transferred to the original algorithm, with a correction term of  $\frac{\delta}{2}$ . ─

### 3.1 Deriving a Simple Probabilistic Invariant

Consider the random variables  $X_s := \mathbf{I}(s \in \chi)$  indicating the presence of a stream element  $s \in A = \{a_1, \dots, a_l\}$  in the buffer, where we write  $\mathbf{I}$  for the indicator of a predicate, so  $\mathbf{I}(\text{true}) = 1$  and  $\mathbf{I}(\text{false}) = 0$ . Before the algorithm first encounters the stream element  $s$ ,  $X_s$  will be 0 unconditionally, because the buffer  $\chi$  is always a subset of the stream elements processed so far, i.e.,  $\chi \subseteq \{a_1, \dots, a_m\}$  after loop iteration  $m$ .

In the loop iteration where element  $s$  occurs for the first time, it will be inserted with probability  $p$  in Lines 3–7. This means, after Line 7, we have:

$$\mathbf{E}[p^{-1}X_s] = 1. \quad (2)$$

Interestingly, this equation is preserved for the rest of the algorithm. For example, let us consider a subsampling step: each  $s$  is independently discarded with probability  $\frac{1}{2}$  so  $\mathcal{P}(X_s = 1)$  is halved, but so is  $p$  after subsampling, which preserves the equation.

Let us see how we can verify Equation 2 more formally. For that, we model the state of the randomized algorithm as a pair  $(\chi, p)$  and we write  $\chi$  and  $p$  for the random variables projecting their respective components from the distribution of the state of the algorithm. We will refer to parts of each loop iteration in Algorithm 1 as  $\text{step}_1$  (resp.  $\text{step}_2$ ) for Lines 3–7 (resp. Lines 8–10). The final distribution of the algorithm is the distribution resulting from the sequential composition of alternating steps over the stream:

$$\text{init} \gg \text{step}_1 a_1 \gg \text{step}_2 \gg \text{step}_1 a_2 \gg \dots \gg \text{step}_1 a_l \gg \text{step}_2$$

where we parameterize  $\text{step}_1$  with the stream element that it processes. The term  $\text{init}$  represents the initial state, i.e.,  $\text{init} = \mathbf{return}(\{\}, 1)$ . It is easy to show by induction over the sequence of steps, we have  $0 < p \leq 1$  and  $\chi \subseteq A$  for all possible states of the algorithm.

Let us verify that Equation 2 is preserved as an invariant over all steps. To verify that  $\text{step}_1 a$  preserves Equation 2, we assume some probability space of states  $\Omega$  fulfills Equation 2 and we would like to show that it is still true for  $\Omega \gg \text{step}_1 a$ . By Equation 1,

$$\mathbf{E}_{\Omega \gg \text{step}_1 a}[p^{-1}X_s] = \int_{\Omega} \int_{\text{Ber}(p)} p^{-1} \mathbf{I}(s \in (\mathbf{if} \ \tau \ \mathbf{then} \ \chi \cup \{a\} \ \mathbf{else} \ \chi - \{a\})) \, d\tau \, d\sigma. \quad (3)$$

Note that we write  $p$  or  $\chi$  even though we should actually write  $p(\sigma)$  or  $\chi(\sigma)$ , i.e., we remember that these implicitly depend on  $\sigma$ . To see that the right-hand side is equal to 1, it is useful to consider cases on whether  $a = s$ . When  $a = s$ , the right-hand-side is equal to 1 by definition of the Bernoulli distribution (since  $p \in (0; 1]$ ). When  $a \neq s$ , it follows from the induction hypothesis on  $\Omega$ ; in particular, the term in the inner integral is constant with respect to  $\tau$ .

The same invariant-based argument is possible for  $\text{step}_2$ . Let us assume  $\Omega$  is a probability space of states fulfilling Equation 2. Then by Equation 1,  $\mathbf{E}_{\Omega \gg \text{step}_2}[p^{-1}X_s]$  equals

$$\int_{\Omega} \left( \mathbf{if} \ |\chi| = n \ \mathbf{then} \ \left( \int_{\text{subsample}(\chi)} \frac{\mathbf{I}(s \in \tau)}{p/2} \, d\tau \right) \ \mathbf{else} \ \frac{\mathbf{I}(s \in \chi)}{p} \right) \, d\sigma. \quad (4)$$

Note that the *true* and *false* cases of the inner **if-then-else** both evaluate to the same value:  $p^{-1} \mathbf{I}(s \in \chi)$ . If  $s \notin \chi$  both sides of the equation are 0, because the subsampling operation returns a subset of  $\chi$ . If  $s \in \chi$  the probability that the element gets subsampled is  $1/2$ , so we arrive again at  $\frac{1/2}{p/2} = p^{-1} \mathbf{I}(s \in \chi)$ . Hence:  $\mathbf{E}_{\Omega \gg \text{step}_2}[p^{-1}X_s] = \mathbf{E}_{\Omega}[p^{-1}X_s] = 1$ . This completes the invariance proof for Equation 2.

### 3.2 Deriving a Functional Probabilistic Invariant

With Equation 2 established, it is straightforward to show that the expected value of the output estimate  $p^{-1}|\chi|$  for the modified algorithm (without Line 11) is equal to the desired cardinality  $|A|$ . However, recall that we are interested in verifying the estimate's PAC guarantee. A typical approach to establishing such a guarantee is to use Chernoff bounds which provide exponential tail bounds (i.e., concentration bounds) for the deviation of sums of independent random variables from their mean. However, these are not directly useful in the CVM algorithm because the key random variables, e.g.,  $p^{-1}X_s$  for  $s \in A$ , are dependent.

An alternative is the Cramér–Chernoff method, which is a general method to obtain tail bounds for any random variable. It can be stated simply as  $\mathcal{P}(X \geq a) \leq M(t)e^{-ta}$  for all  $t > 0$ , where  $M(t) := \mathbb{E}[\exp(tX)]$  is the moment generating function of the random variable  $X$ . It is also possible to obtain lower tail bounds  $\mathcal{P}(X \leq a)$  using the Cramér–Chernoff method, which just requires estimates for  $M(t)$  for  $t < 0$ , instead of  $t > 0$ .

In our case, we are interested in estimating the moment generating function of the random variable  $p^{-1}|\chi|$  for the CVM algorithm:

$$\mathbb{E}[\exp(tp^{-1}|\chi|)] = \mathbb{E}\left[\prod_{s \in A} h(p^{-1}X_s)\right]$$

for  $h(x) = \exp(tx)$ . At this point, it is tempting to see whether the proof for Equation 2 can be generalized to establish bounds for the above. Indeed, we managed to establish the following generalized result:

$$\mathbb{E}\left[\prod_{s \in A} h(p^{-1}X_s)\right] \leq h(1)^{|A|} \quad (5)$$

for every non-negative concave function  $h : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ . However, the exponential function in  $M(t)$  is convex, but we can instead try to derive tail bounds for the random variable  $I(p \geq q)p^{-1}|\chi|$ , for some fixed constant  $q > 0$ . This leads to a similar invariant inequality:

$$\mathbb{E}\left[\prod_{s \in A} I(p \geq q)h(p^{-1}X_s)\right] \leq h(1)^{|A|} \quad (6)$$

with the new condition that  $h$  needs to be non-negative and concave only on  $[0; q^{-1}]$ . This then allows us to approximate the exponential function from above with an affine function  $h$  on the range  $[0; q^{-1}]$ , which yields tail bounds for  $p^{-1}|\chi|$  under the condition  $p \geq q$ . As an example, the upper tail bound can be derived as follows:

$$\begin{aligned} \mathcal{P}(p^{-1}|\chi| \geq (1 + \varepsilon)|A| \wedge p \geq q) &\leq \mathcal{P}(I(p \geq q)p^{-1}|\chi| \geq (1 + \varepsilon)|A|) \\ &\stackrel{\text{Markov}}{\leq} e^{-t(1+\varepsilon)|A|} \mathbb{E}\left[\prod_{s \in A} I(p \geq q) \exp(tp^{-1}X_s)\right] \\ &\leq e^{-t(1+\varepsilon)|A|} \mathbb{E}\left[\prod_{s \in A} I(p \geq q)h(p^{-1}X_s)\right] \\ &\stackrel{\text{Ineq. 6}}{\leq} e^{-t(1+\varepsilon)|A|} h(1)^{|A|} \\ &\stackrel{\text{Calculus}}{\leq} e^{-n\varepsilon^2/12} \end{aligned}$$

■ **Algorithm 3** New total and unbiased CVM algorithm variant.

**Input:** Stream elements  $a_1, \dots, a_l$ ,  $0 < \varepsilon$ ,  $0 < \delta < 1$ .

**Output:** A cardinality estimate  $R$  for set  $A = \{a_1, \dots, a_l\}$  such that  $\mathcal{P}(|R - |A|| > \varepsilon|A|) \leq \delta$

```

1:  $\chi \leftarrow \{\}, p \leftarrow 1, n = \lceil \frac{12}{\varepsilon^2} \ln(\frac{3l}{\delta}) \rceil, \frac{1}{2} \leq f < 1$ , such that  $nf$  integer
2: for  $i \leftarrow 1$  to  $l$  do
3:    $b \xleftarrow{\$} \text{Ber}(p)$  ▷ insert  $a_i$  with probability  $p$  (and remove it otherwise)
4:   if  $b$  then
5:      $\chi \leftarrow \chi \cup \{a_i\}$ 
6:   else
7:      $\chi \leftarrow \chi - \{a_i\}$ 
8:   if  $|\chi| = n$  then ▷ if buffer  $\chi$  is full
9:      $\chi \xleftarrow{\$} \text{subsample}(\chi)$  ▷ select a random  $nf$ -subset of  $\chi$ 
10:     $p \leftarrow pf$ 
11: return  $\frac{|\chi|}{p}$  ▷ estimate cardinality of  $A$ 

```

where we choose  $h(x) = 1 + qx(e^{t/q} - 1)$ . Note that  $h$  is affine and it can be easily checked<sup>3</sup> that it is an upper approximation of  $\exp(tx)$  for  $x \in [0; q^{-1}]$ . For the last step, we have to find the  $t$  that produces the required bound.<sup>4</sup> To use these bounds, we also have to separately estimate  $\mathcal{P}(p < q)$ . For that, we use a similar strategy, as in the original proof by Chakraborty et al. [11], with  $q = \frac{n}{4|A|}$ . The formalization in the supplementary material [34] contains a detailed informal step-by-step proof using our approach in its appendix. Besides the use of Equation 1 and the Cramér–Chernoff method, the steps are elementary.

We call inequalities like Inequality 5 and 6: *functional probabilistic invariants*. The inequalities can be established using induction over the steps of the randomized algorithm. Even though the actual distribution of the states themselves are not known, nor do we think it is possible to find useful closed form descriptions for them, the invariant establishes valuable information about the distribution of the state. In this case, enough information, to establish exponential tail bounds for the algorithm.

## 4 An Unbiased CVM Variant and Negative Dependence

An interesting consequence of our invariant-based approach is that it allowed us to devise and verify a refined version of the CVM algorithm that is both total and unbiased.

### 4.1 Unbiased CVM Variant

When we look at the subsampling step of Algorithm 1, our invariant (Inequality 5) imposes the following condition on the subsampling operation. It should be noted that the condition becomes apparent while establishing Inequality 5 using similar but more general steps as described in Equations 3 and 4.<sup>5</sup>

$$\int_{\text{subsample}(\chi)} \prod_{s \in S} g(\mathbf{I}(s \in \tau)) \, d\tau \leq \prod_{s \in S} \mathbb{E}_{\text{Ber}(f)}[g] \quad (7)$$

<sup>3</sup> Because the exponential function is convex and  $h$  is affine, we only have to check the end points:  $0, q^{-1}$ .

<sup>4</sup> We use  $t = q \ln(1 + \varepsilon)$  which is not the real optimum, but better for algebraic evaluation.

<sup>5</sup> A step-by-step proof of the derivation is also available in the appendix of the formalization [34, Appendix A].



for all non-negative functions  $g$  and any  $S \subseteq \chi$ , where  $f$  is the probability of retaining each element in the subsampling step of the algorithm. (This parameter  $f$  was fixed to  $\frac{1}{2}$  in the original presentation of the algorithm for simplicity.) Any subsampling step that satisfies Inequality 7 can be used while still preserving Inequalities 5 and 6 for the algorithm.

Motivated by this observation, our new variant is shown in Algorithm 3. For the subsampling step, instead of keeping each element of  $\chi$  with probability  $\frac{1}{2}$ , we pick a uniform random  $nf$ -subset of  $\chi$ , where  $\frac{1}{2} \leq f < 1$  and  $nf$  is an integer. For example, it is possible to choose  $f = \frac{n-1}{n}$ , i.e., discarding just one random element from  $\chi$  in the subsampling step. Since this new subsampling step always reduces the size of  $\chi$ , the variant is *total* (never returns  $\perp$ ). The invariant-based approach allows us to show that the algorithm is probably-approximately correct and also *unbiased*, i.e., the expectation of the result is exactly  $|A|$ . These depend crucially on establishing Inequality 7 for the new subsampler, for which we need a new concept.

## 4.2 Background on Negative Dependence

Some sets of random variables possess a property called *negative association*, a generalization of independence. The concept was introduced by Joag-Dev and Proschan [29], who showed that it has many useful closure properties compared to other previously introduced notions of negative dependence, such as negative correlation or negative orthant dependence. Importantly, standard Chernoff–Hoeffding type bounds still apply to negatively associated random variables [16, Prop. 7]. Negative association is defined as follows:

► **Definition 3.** For a function defined on  $n$ -tuples  $f : V^n \rightarrow W$ , we will denote by  $\text{dep}(f)$  the set of coordinates the function depends on, i.e.,  $\text{dep}(f) \subseteq \{1, \dots, n\}$  is minimal, such that  $f(x) = f(y)$  for all  $x, y \in V^n$  with  $x_i = y_i$  for all  $i \in \text{dep}(f)$ .

► **Definition 4 (Negative Association).** A set of random variables  $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$  is negatively associated if, for all non-decreasing functions  $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ , which depend on disjoint sets of the variables, i.e.,  $\text{dep}(f) \cap \text{dep}(g) = \emptyset$ , the following inequality holds:

$$\mathbb{E}[f(X_1, \dots, X_n)g(X_1, \dots, X_n)] \leq \mathbb{E}[f(X_1, \dots, X_n)]\mathbb{E}[g(X_1, \dots, X_n)].$$

The following proposition summarizes some important properties of negatively associated sets of random variables.

► **Proposition 5** (Summary of results for negatively associated random variables [29]).

1. If  $X = (X_1, \dots, X_n)$  are negatively associated then  $\mathbb{E}[f(X)g(X)] \leq \mathbb{E}[f(X)]\mathbb{E}[g(X)]$  for non-increasing functions  $f, g$  with  $\text{dep}(f) \cap \text{dep}(g) = \emptyset$ .
2. If  $X = (X_1, \dots, X_n)$  are negatively associated,  $Y = (Y_1, \dots, Y_m)$  are negatively associated, and the pair of vector-valued random variables  $X$  and  $Y$  are independent, then the union  $X_1, \dots, X_n, Y_1, \dots, Y_m$  is a set of negatively associated random variables.
3. If  $X = (X_1, \dots, X_n)$  are negatively associated and  $f_1, \dots, f_m : \mathbb{R}^n \rightarrow \mathbb{R}$  are all non-increasing or all non-decreasing functions, such that  $\text{dep}(f_i) \cap \text{dep}(f_j) = \emptyset$  for  $i \neq j$ , then  $f_1(X), \dots, f_m(X)$  form a set of negatively dependent random variables of size  $m$ .
4. If  $X_1, \dots, X_n$  are independent then  $X_1, \dots, X_n$  are negatively associated.
5. A subset of a negatively associated set of random variables is again negatively associated.

These properties illustrate the trade-off between negative association and independence. For example, Property 3 would be true for independent random variables, even without the condition of monotonicity. To analyze our new subsampler, the following is an important lemma about negatively associated random variables.

► **Lemma 6.** *Let  $X_1, \dots, X_n$  be negatively associated and  $f_1, \dots, f_n$  be all non-decreasing (or all non-increasing), non-negative functions, then*

$$\mathbb{E} \left[ \prod_{i=1}^n f_i(X_i) \right] \leq \prod_{i=1}^n f_i(\mathbb{E}[X_i]).$$

**Proof.** This follows from the definition of negative association (or Property 1 of Proposition 5, if the  $f_i$  are non-increasing) using induction. ◀

The case for non-decreasing functions of the above lemma is pointed out by Joag-Dev and Proschan [29, P.2]. The reason for our interest in this lemma stems from the fact that indicator variables of random  $m$ -subsets are negatively associated. This is a consequence of the fact that permutation distributions are negatively associated [29, Th. 2.11]. Thus, for the new subsampling step in Line 9 of Algorithm 3, we can derive using Lemma 6:

$$\int_{\text{subsample}(\chi)} \prod_{s \in S} g(\mathbb{I}(s \in \tau)) \, d\tau \leq \prod_{s \in S} \int_{\text{subsample}(\chi)} g(\mathbb{I}(s \in \tau)) \, d\tau = \prod_{s \in S} \mathbb{E}_{\text{Ber}(f)}[g]. \quad (8)$$

for any non-negative  $g$  and  $S \subseteq \chi$ . Note that the domain of  $g$  has two values, so it is either non-increasing or non-decreasing. Also, if  $S$  is a singleton, the inequality becomes an equality. With this ingredient, we can conclude that our results about the original algorithm derived in the previous section also hold for our new variant (Algorithm 3).

## 5 Formalization of the CVM algorithm

Let us now turn to details of our formalization of the CVM algorithm in Isabelle/HOL using our invariant-based approach [34]. We verified both the total, unbiased variant (Algorithm 3) and the original variant (Algorithm 1) from the introduction.

► **Note 7.** In our supplementary material [34], the theory `CVM_Abstract_Algorithm` verifies a generalized version of the CVM algorithm, with an abstract subsampling operation that is required to fulfill Inequality 7. The specialization happens in the following theories, where `CVM_Original_Algorithm` verifies the original algorithm, and `CVM_New_Unbiased_Algorithm` verifies the new total and unbiased variant. Note that only `CVM_New_Unbiased_Algorithm` depends on the new library for negatively associated random variables, which we describe in more detail in Section 6. The total number of lines required for the verification of the original algorithm is 1003 lines. In addition, we actually verified a slight generalization of Algorithm 1 where the subsampling probability can be any  $f \in [\frac{1}{2}; e^{-1/12}]$ ; the original CVM algorithm [10] is the special case  $f = \frac{1}{2}$ . ◀

A snippet of the formalization of Algorithm 3 is presented in Figure 1 (the formalization of Algorithm 1 is very similar). We use the same variables as in the informal presentation:  $n$  for the maximal size of the buffer,  $f$  for the fraction of elements to keep in the buffer when subsampling. The condition  $\langle n * f \in \mathbb{N} \rangle$  expresses the requirement that the  $nf$  must be integer. Instead of representing the state using pairs, as we did in the informal discussion, we use a datatype with the single constructor `State`, which has two arguments  $\chi$  and  $p$ , the buffer and the probability that the stream elements are in the buffer, respectively. Isabelle/HOL provides notation closely related to informal pseudocode, so it is usually feasible to read a formal statement without expert knowledge. Nevertheless, Table 1 contains a brief glossary of the syntax used in the formalization.

The theorem that establishes the correctness of the algorithm, i.e., that the relative error will exceed  $\varepsilon$  with probability at most  $\delta$  is expressed in the following snippet:

```

context
  fixes  $f :: \text{real}$  and  $n :: \text{nat}$ 
  assumes  $f\text{-range}$ :  $\langle f \in \{1/2..<1\} \rangle \langle n * f \in \mathbb{N} \rangle$  and  $n\text{-gt-0}$ :  $\langle n > 0 \rangle$ 
begin

  definition  $\langle \text{initial-state} = \text{State } \{\} \ 1 \rangle$  — Setup initial state  $\chi = \emptyset$  and  $p = 1$ .
  fun subsample where — Subsampling operation: Sample random  $nf$  subset.
     $\langle \text{subsample } \chi = \text{pmf-of-set } \{S. S \subseteq \chi \wedge \text{card } S = n * f\} \rangle$ 

  fun step where — Loop body.
     $\langle \text{step } a \ (\text{State } \chi \ p) = \text{do } \{$ 
       $b \leftarrow \text{bernoulli-pmf } p;$ 
       $\text{let } \chi = (\text{if } b \text{ then } \chi \cup \{a\} \text{ else } \chi - \{a\});$ 

       $\text{if } \text{card } \chi = n \text{ then } \text{do } \{$ 
         $\chi \leftarrow \text{subsample } \chi;$ 
         $\text{return-pmf } (\text{State } \chi \ (p * f))$ 
       $\} \text{ else } \text{do } \{$ 
         $\text{return-pmf } (\text{State } \chi \ p)$ 
       $\}$ 
     $\rangle$ 

  fun run-steps where — Iterate loop over stream  $xs$ .
     $\langle \text{run-steps } xs = \text{foldM-pmf step } xs \ \text{initial-state} \rangle$ 
  fun estimate where  $\langle \text{estimate } (\text{State } \chi \ p) = \text{card } \chi / p \rangle$ 
  fun run-algo where — Run algorithm and estimate.
     $\langle \text{run-algo } xs = \text{map-pmf estimate } (\text{run-steps } xs) \rangle$ 
  [...]
end

```

■ **Figure 1** Formalized version of Algorithm 3.

```

theorem correctness:
  assumes  $\langle \varepsilon \in \{0 < .. < 1 :: \text{real}\} \rangle \langle \delta \in \{0 < .. < 1 :: \text{real}\} \rangle$ 
  assumes  $\langle \text{real } n \geq 12 / \varepsilon^2 * \ln (3 * \text{real } (\text{length } xs) / \delta) \rangle$ 
  defines  $\langle A \equiv \text{real } (\text{card } (\text{set } xs)) \rangle$ 
  shows  $\langle \mathcal{P}(R \text{ in run-algo } xs. |R - A| > \varepsilon * A) \leq \delta \rangle$ 

```

The first line gives conditions on parameters  $\varepsilon$  and  $\delta$ , which must be strictly between 0 and 1. The next line requires the buffer size  $n$  to be larger than or equal to  $12\varepsilon^{-2} \ln(3\delta^{-1}l)$ . Then, we introduce the abbreviation  $A$  for the cardinality of the set of elements in the sequence  $xs$ . The notation  $\mathcal{P}(x \text{ in } M. P \ x)$  denotes the probability of a predicate  $P$  in the probability space  $M$ , so the conclusion gives the PAC guarantee for the output estimate  $R$  from *run-algo*.

Similarly, we have also formalized unbiasedness of Algorithm 3:

```

theorem unbiasedness:  $\langle \text{measure-pmf.expectation } (\text{run-algo } xs) \ id = \text{card } (\text{set } xs) \rangle$ 

```

where the expression *measure-pmf.expectation*  $M \ f$  denotes the expectation of the random variable  $f$  on the probability space  $M$ .

Our proofs are available both in mechanized form in Isabelle/HOL and as a pen-and-paper proof included in the associated proof document. In practice, we developed the latter proof first and then mechanized it in Isabelle/HOL without much surprise. Most of the lemmas can be identified one-to-one between both proofs; Isabelle’s existing libraries, automation capabilities, and structured proof format were used extensively in our proofs.

■ **Table 1** Isabelle/HOL syntax used in Figure 1.

Term	Description
<i>card</i> $S$	Cardinality of a finite set $S$ .
<i>real</i>	Type of real numbers and conversion from natural numbers into real numbers.
<i>nat</i>	Type of natural numbers (non-negative integers).
<i>bernoulli-pmf</i> $p$	The probability space over the Boolean values, where the probability of <i>True</i> is $p$ . (Bernoulli distribution.)
<i>pmf-of-set</i> $S$	For a finite set $S$ , the uniform probability space on $S$ . (Every element of $S$ is equiprobable.)
<i>map-pmf</i> $f$ $A$	The probability space representing the distribution of the random variable $f$ over the probability space $A$ .
<i>return-pmf</i> $x$	The probability space of the singleton $\{x\}$ .
<i>foldM-pmf</i> $f$ $xs$ $a$	Iterate randomized algorithm $f$ over the sequence $xs$ using the initial state $a$ .

## 6 Formalization of a Library for Negative Association

As mentioned in Section 4, formalizing the total and unbiased variant of the CVM algorithm requires results from the theory of negative association.

► **Note 8.** The formalization of the theory of negative association is included in a separate AFP entry [33]. This library contains key results used to establish the invariants for CVM (e.g., `Neg_Assoc_Permutation_Distributions`). Although not needed for CVM, we have also mechanized the standard Chernoff bounds (`Neg_Assoc_Chernoff_Bounds`), including the additive bounds by Hoeffding [27, Th. 1, 2] and the multiplicative bounds by Motwani and Raghavan [40, Th. 4.1, 4.2]. Another example application included in the library is proving the false positive rate of Bloom filters (`Neg_Assoc_Bloom_Filters`). In total, the library contains 2974 lines of Isabelle code.  $\lrcorner$

Our formalization follows the definitions by Joag-Dev and Proschan [29] closely. However, their definition leaves the class of test functions  $f$  and  $g$  (in Definition 4) imprecise. In the formalization, we use different conditions for introduction and elimination rules. In particular, for introduction rules, the test functions are bounded and measurable. However, we provide stronger elimination rules, showing that if  $X_1, \dots, X_n$  are negatively associated, then the inequality on expectations is true even if  $f, g$  are only square-integrable; or, alternatively, integrable and non-negative. This is derived using the monotone convergence theorem.

Another deviation from the original work is that we do not require that the random variables are real-valued. In the formalization, any linearly ordered topological space with the Borel  $\sigma$ -algebra is allowed as the range space. In this case, the test functions must be monotone with respect to the respective order on the range space.

A key issue we faced during formalization was that there are many theorems that condition on a set of functions being either simultaneously monotone or simultaneously anti-monotone. To reduce duplication, we introduce a notation that allows us to abstract over the direction of relations:  $\leq_{\geq \eta}$ ; it evaluates to the forward version of the relation  $\leq$  if  $\eta = \text{Fwd}$  and the converse:  $\geq$  if  $\eta = \text{Rev}$ . For example the FKG inequality [3, Ch. 6],[20]

$$\mathbb{E}[fg] \geq \mathbb{E}[f] \mathbb{E}[g]$$

is true, if  $f$  and  $g$  are both monotone, or both antimonotone, on a probability space whose domain is a finite distributive lattice with a log-supermodular probability mass function. The reverse inequality is also true, if  $f$  is monotone and  $g$  is antimonotone, or vice versa. Using our parameterized relation symbol, we can state all variants in a concise manner:

```
theorem fkf-inequality-pmf:
  fixes  $M :: \langle 'a :: \text{finite-distrib-lattice} \rangle \text{ pmf} \rangle$ 
  fixes  $f\ g :: \langle 'a \Rightarrow \text{real} \rangle$ 
  assumes  $\langle \bigwedge x\ y. \text{pmf } M\ x * \text{pmf } M\ y \leq \text{pmf } M\ (x \sqcup y) * \text{pmf } M\ (x \sqcap y) \rangle$ 
  assumes  $\langle \text{monotone } (\leq) (\leq_{\geq \tau}) f \rangle \langle \text{monotone } (\leq) (\leq_{\geq \sigma}) g \rangle$ 
  shows  $\langle (\int x. f\ x\ \partial M) * (\int x. g\ x\ \partial M) \leq_{\geq \tau * \sigma} (\int x. f\ x * g\ x\ \partial M) \rangle$ 
```

Here,  $\sigma$  and  $\tau$  are relation directions, and  $\sigma * \tau$  multiplies relation directions, i.e.,  $\sigma * \tau$  is the forward direction if  $\sigma$  and  $\tau$  have the same direction, and it is the reverse direction otherwise. The first assumption denotes the log-supermodularity of the probability mass function, while the second assumptions are the parametric monotonicity conditions. The FKG inequality is a key result which enables verification of negative association for random variables. This includes the indicator variables for the new subsampling operation we introduced in Section 4.

Let us summarize a few key formalized results for negatively associated random variables in our library. The following is the well-known Hoeffding inequality [27] generalized for negatively associated random variables.

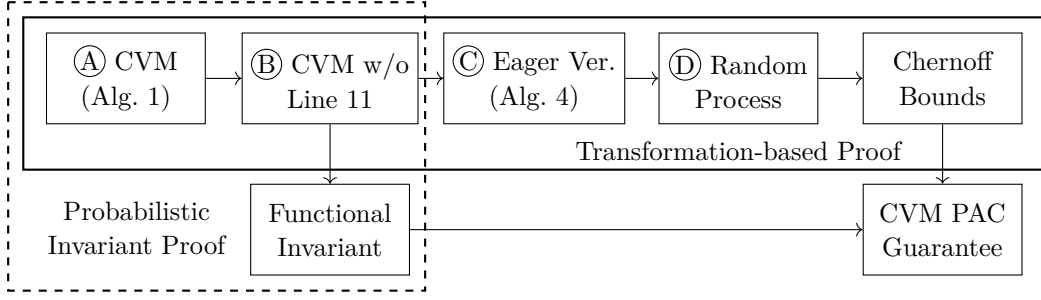
```
lemma hoeffding-bound-two-sided:
  assumes  $\langle \text{neg-assoc } X\ I \rangle \langle \text{finite } I \rangle$ 
  assumes  $\langle \bigwedge i. i \in I \Rightarrow a\ i \leq b\ i \rangle$ 
  assumes  $\langle \bigwedge i. i \in I \Rightarrow AE\ \omega\ \text{in } M. X\ i\ \omega \in \{a\ i..b\ i\} \rangle \langle I \neq \{\} \rangle$ 
  defines  $\langle n \equiv \text{real } (\text{card } I) \rangle$ 
  defines  $\langle \mu \equiv (\sum i \in I. \text{expectation } (X\ i)) \rangle$ 
  assumes  $\langle \delta \geq 0 \rangle \langle (\sum i \in I. (b\ i - a\ i)^2) > 0 \rangle$ 
  shows  $\langle \mathcal{P}(\omega\ \text{in } M. |\sum i \in I. X\ i\ \omega - \mu| \geq \delta * n) \leq 2 * \exp(-2 * (n * \delta)^2 / (\sum i \in I. (b\ i - a\ i)^2)) \rangle$ 
```

Another key result (shown below) used for the verification of our CVM variant is the negative-associativity of the indicator functions of random  $k$ -subsets of a finite set  $S$  (with cardinality greater than or equal to  $k$ ).

```
lemma n-subsets-distribution-neg-assoc:
  assumes  $\langle \text{finite } S \rangle \langle k \leq \text{card } S \rangle$ 
  defines  $\langle p \equiv \text{pmf-of-set } \{T. T \subseteq S \wedge \text{card } T = k\} \rangle$ 
  shows  $\langle \text{measure-pmf.neg-assoc } p\ (\in)\ S \rangle$ 
```

This is a consequence of a more general result, which we have also shown, that permutation distributions are negatively associated. We relied on the proof by Dubhashi et al. using the FKG inequality [15, Th. 10]; there is a prior proof by Joag-Dev and Proschan [29, Th. 2.11], which is incomplete.<sup>6</sup>

<sup>6</sup> The step which we could not directly formalize is the assertion (Sentence 14) in the proof of Theorem 2.11 ([29]) that the conditional expectation of the random variable  $f(X)$  is smaller iff the smallest element of the permutation is contained in  $\text{dep}(f)$ . That statement is non-trivial and requires a proof using a theorem such as the FKG-inequality. We think Dubhashi et al. developed their proof to complete it.



■ **Figure 2** An overview of the two formalized proof approaches for the CVM algorithm.

## 7 Transformation-Based Proof

Here, we describe the transformation-based proof by Chakraborty et al. [11], focusing on the challenging parts in its formalization. An overview of the proof is shown in Figure 2, which highlights, in part, why the transformation-based approach required more work to formalize, and why we developed our new approach using probabilistic invariants.

► **Note 9.** The transformation-based formalization of the CVM algorithm is included in the supplementary material [48]. To formalize probabilistic transformations (relating two distributions), we built on an existing relational program logic in Isabelle/HOL [37]. The formalization took 2634 lines which is considerably longer than the proof using our invariant-based technique (1003 lines). ┘

As mentioned in Section 1, the main difficulty in directly analyzing Algorithm 1 is the lack of independence in its state variables. The technique Chakraborty et al. use to circumvent this issue is by progressively modifying the algorithm ((A)–(D) in Figure 2), in a manner that obviously bounds (or preserves) its distribution, and such that the final algorithm (D) can be described using a simple random process with independent coin flips.

For interested readers, (A) corresponds to [11, Algorithm 1], (B) is [11, Algorithm 2], and (D) is [11, Algorithm 3]. Whereas Chakraborty et al. move directly from (B) to (D) with an informal argument, (C) is a transformation we added in the formalization to bridge this gap.

### 7.1 A Bridging Transformation

Let us consider algorithm (B) in a state where  $k$  subsampling steps have been performed, i.e.,  $p = 2^{-k}$ . The algorithm would perform a coin flip lazily with probability  $p$  when it encounters the next stream element. The transformation (C) is shown in Algorithm 4, and we prove that it computes precisely the same distribution as (B). In (C), we eagerly perform a fixed number of coin flips for each sequence element at the beginning. Now, each element is put into the state  $\chi$ , whenever the first  $k$  coin flips associated with the sequence element are all 1s. This happens exactly with probability  $2^{-k}$ , which means the behaviour of the algorithm is unchanged from (B). Similarly, in the subsampling operation, only those elements whose  $k + 1$ -th associated coin flip is 1 are kept; the operation  $p \mapsto \frac{p}{2}$  is replaced with  $k \mapsto k + 1$ . This again preserves the behaviour of (B) that each element is discarded independently with probability  $1/2$ .

It is easy to show for (C) that the coin flips are independent, and that the set of elements in  $\chi$  in any state are exactly those stream elements for which the first  $k$  entries of their associated coin flips are 1. The final random process (D) directly computes the final set of elements in  $\chi$  after the stream, taking  $K$  as a fixed parameter; one relates (C) to (D) by:

$$\mathcal{P}_{\textcircled{C}}(k = K \wedge \chi = X) \leq \mathcal{P}_{\textcircled{D}_K}(\chi = X)$$

■ **Algorithm 4** Modified CVM algorithm with independent coin flips. The function `last_index` returns the index of the last occurrence of an element in the sequence, before the current iteration.

**Input:** Stream elements  $a_1, \dots, a_l$ ,  $0 < \varepsilon$ ,  $0 < \delta < 1$ .

**Output:** A cardinality estimate  $R$  for set  $A = \{a_1, \dots, a_l\}$  such that  $\mathcal{P}(|R - |A|| > \varepsilon|A|) \leq \delta$

```

1:  $\chi \leftarrow \{\}, k \leftarrow 0, n = \lceil \frac{12}{\varepsilon^2} \ln(\frac{6l}{\delta}) \rceil$ 
2:  $b[i, j] \xleftarrow{\$} \text{Ber}(1/2)$  for  $i, j \in \{1, \dots, l\}$   $\triangleright$  perform  $l^2$  unbiased independent coin flips
3: for  $i \leftarrow 1$  to  $l$  do
4:   if  $b[i, 1] = b[i, 2] = \dots = b[i, k] = 1$  then  $\triangleright$  insert  $a_i$  if first  $k$  flips are 1s.
5:      $\chi \leftarrow \chi \cup \{a_i\}$ 
6:   else
7:      $\chi \leftarrow \chi - \{a_i\}$ 
8:   if  $|\chi| = n$  then  $\triangleright$  if buffer  $\chi$  is full
9:      $\chi \leftarrow \{a \in \chi \mid b[\text{last\_index}(a), k+1] = 1\}$   $\triangleright$  keep elems. whose  $k+1$ -th flip is 1
10:     $k \leftarrow k + 1$ 
11: return  $2^k |\chi|$   $\triangleright$  estimate cardinality of  $A$ 
```

for fixed values of  $K$  and  $X$ . To see how tail bounds can be derived from this inequality, let us first consider the failure event where the algorithm  $\textcircled{C}$ 's estimate exceeds the desired estimation interval and it ends with some fixed value  $k = K$ . Using  $\textcircled{D}$ , this can be bounded using a Chernoff bound for the probability that the number of stream elements whose associated coin flips start with  $K$  1s is outside  $2^{-K}|A|(1 \pm \varepsilon)$ . Now, we can take a union bound over all the possible values  $K$  to establish a global bound for the failure event in  $\textcircled{C}$ . This is explained in more detail by Chakraborty et al. [11].

## 7.2 Eager to Lazy Coin Flips

A remaining question is how to formalize the transformation from  $\textcircled{B}$  to  $\textcircled{C}$ . Our insight is that it is best to solve the problem backwards, i.e., we start with the modified algorithm  $\textcircled{C}$ , which performs all the coin flips in advance *eagerly* and convert it back to  $\textcircled{B}$  which implicitly performs the coin flips *lazily* at the point they are needed.

The main idea is to automatically push down the coin flips through the expression tree of Algorithm 4. To explain how this works, let us first define the *sampling* function, i.e., let  $f$  be a function that takes as argument a vector of coin flips indexed by  $I$ , then we can express the distribution of  $f$  with respect to independent unbiased coin flips as:

$$\text{sample } f = \text{map-pmf } f \text{ (prod-pmf } I \text{ (}\lambda\cdot. \text{bernoulli-pmf } (1/2)\text{))}$$

The interesting fact is that we can distribute the sampling operation over composition:

► **Observation 10.** Let  $f, g$  be functions consuming a set of coin flips (indexed by  $I$ ), where  $g$  also consumes the output of  $f$ , such that,  $f$  depends only on the coin flips indexed by  $J \subseteq I$  and  $g$  depends on the complement  $I - J$ , then:

$$\text{sample } (\lambda\omega. g \omega \circ f \omega) = \text{sample } f \gg (\lambda x. \text{sample } (\lambda\omega. g \omega x))$$

By recursively applying the observation, we end up with elementary lookup operations, e.g.,  $\text{sample } (\lambda\omega. \omega i)$ , for which it is easy to see that it is just a coin flip, i.e., equal to  $\text{bernoulli-pmf } (1/2)$ . This lets us readily transform  $\textcircled{C}$  to  $\textcircled{B}$  and prove their distributions equivalent.



A detail that we have simplified here is that the split of the index sets, e.g., which coin flips  $f$  depends on and which coin flips  $g$  depends on, may be dynamic. For example, when the algorithm increases the subsampling counter  $k$ , it will have read the corresponding row of coin flips. This means we have a situation where the previous loop iteration communicates to the next loop iteration which coin flips it depends on using the state; and the next loop iteration will indeed only read coin flips that were not read by the previous iteration.

To handle these situations we generalized Observation 10 to allow for the case where the set of indices  $J$  splitting the set of coin flips  $f$  and  $g$  depend on, may itself depend on the result of  $f$ .

## 8 Related Work

### 8.1 Algorithms for the Distinct Elements Problem

It is important to note that there are several practical solutions for the distinct elements problem. The first solution was presented by Flajolet [19] in 1985; however, like many other authors [18, 26], his solution makes the assumption that a fixed hash function can be regarded as a fully random function. Alon et al. [2, Section 2.3] presented an easy remedy, which does not require such unmotivated model assumptions. Their algorithm just relies on keeping track of the maximum of the hash values of the stream elements, where the hash function must be chosen uniformly from a pairwise independent family.

Later, Bar-Yossef et al. [5], Kane et al. [30] and Błasiok in 2020 [8] improved on the solution by Alon et al. For example, Bar-Yossef et al. present a solution (Algorithm 3 in their work) with a space-complexity of  $\mathcal{O}(\ln(\delta^{-1})(\varepsilon^{-2}(\ln(\varepsilon^{-1}) + \ln b) + b))$ , which can be implemented in practice. This is slightly better than the CVM algorithm which requires  $\mathcal{O}(\varepsilon^{-2} \ln(\delta^{-1}l)b)$ . In particular, there is no dependency on the length of the stream  $l$ . The more recent and more sophisticated solution by Błasiok is space-optimal, with a space complexity of  $\mathcal{O}(\varepsilon^{-2} \ln(\delta^{-1}) + b)$ . We [32] presented a version of the latter that preserves monotonicity and supports the merge-operation, which enables its use in distributed settings, such as Map-Reduce pipelines [13]. It should be noted that these recent algorithms are mostly of theoretical interest, as the constants, as well as the implementation complexity, are rather large. A comprehensive review of distinct elements algorithms has been compiled by Pettie and Wang [43, Table 1]. What makes the CVM algorithm unique is its simplicity and the fact that it does not rely on hashing, which may enable more general use-cases than the traditional algorithms.

The aforementioned hash-based algorithms are biased; Flajolet et al. [19] points this out and also provides bounds on the distance between the expected result and the cardinality of the stream. Most authors do not discuss the matter of bias but it is not hard to show. One issue, for example, is that the usual method to amplify the accuracy of these algorithms is using the median, which does not preserve expectations. In the context of query processing, unbiasedness has been discussed [23, Section 2.1], but we could not find any similar discussion for the distinct elements problem in the streaming model.

### 8.2 Probabilistic Invariants and Formalization

As far as we know, probabilistic invariants have not been used to establish exponentially decreasing tail-bounds. However, it is fairly common to use recursive analysis techniques to establish results about expectations or variance of random variables, such as their run-time [40, Section 1.4]. This is easy due to the linearity of expectations and – for independent random variables – linearity of variances. A simple example is the Morris counter [44] or the expected run-time of the quick-sort algorithm [39, Section 2.5].



There is also research on the (automated) analysis of loop invariants, for probabilistic loops, using their characteristic functions [7, 38]. This approach works by establishing the limiting distribution of the state of the loop. De Medeiros et al. [12, Section 3.2] also establish methods to derive limiting distributions of probabilistic loops. Our approach differs from these techniques by avoiding computation of the distribution, which, we think, is infeasible for the CVM algorithm. Instead, we investigate invariants of classes of functions of the distributions, which are relevant for the analysis. There is research on automated evaluation of moments for restricted classes of loops which contain only polynomial assignments and no branches [6, 36]. However, these methods do not extend to algorithms with non-continuous operations, or control flow that depends on non-deterministic state variables.

Finally, verification of randomized algorithms expressed in functional programming languages has been tackled by various authors using various proof assistants [4, 9, 12, 17, 22, 25, 28, 45, 46, 47]; the most closely related efforts are our mechanizations of frequency moments algorithms [31, 32]. Moreover, for classical imperative randomized algorithms, there are approaches based on probabilistic Hoare logic [14]. An interesting related work by Haselwarter et al. [24] discusses some of the issues we encountered in the transformation-based proof (Section 7), such as (approximate) equivalence between related randomized algorithms. However, because our work reasons directly over the semantics of randomized algorithms expressed in the Giry monad, we did not deeply explore probabilistic logic-based methods.

## 9 Conclusion

We presented the first formalization of the CVM algorithm using Isabelle/HOL. Central to our formalization is a novel invariant-based proof technique to establish exponentially decreasing tail-bounds for randomized algorithms, which is inspired by our alternative analysis of the CVM algorithm via the Cramér–Chernoff method. Our technique can be summarized by the following two steps:

1. Find functionals over the state distribution of the algorithm, for which it is possible to establish bounds on their expectation inductively/recursively.
2. Use those bounds to establish tail bounds on the result of the algorithm.

Comparing our approach against the original proof by Chakraborty et al. [11] shows that our technique yields a considerably shorter formalization (with 1003 vs. 2634 lines). Interestingly, our technique also readily generalized to a new CVM variant with stronger properties (totality and unbiasedness) – we formalized this latter version using the same invariant, together with a new library of results for negative association. In future work, it would be interesting to formalize other variations of subsampling for CVM.

Note that we have yet to apply our proof technique to examples beyond CVM. (It is easy to construct artificial examples.) Identifying realistic applications for our new method is an interesting avenue for future work – this could lead to further refinements of the method, and a better understanding of how to identify suitable functionals for the proofs.

---

## References

- 1 Archive of Formal Proofs. Accessed: 2025-02-26. URL: <https://isa-afp.org>.
- 2 Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. *J. Comput. Syst. Sci.*, 58(1):137–147, 1999. doi:10.1006/JCSS.1997.1545.
- 3 Noga Alon and Joel H. Spencer. *The Probabilistic Method, Second Edition*. John Wiley, 2000. doi:10.1002/0471722154.

- 4 Philippe Audebaud and Christine Paulin-Mohring. Proofs of randomized algorithms in Coq. *Sci. Comput. Program.*, 74(8):568–589, 2009. doi:10.1016/J.SCICO.2007.09.002.
- 5 Ziv Bar-Yossef, T. S. Jayram, Ravi Kumar, D. Sivakumar, and Luca Trevisan. Counting distinct elements in a data stream. In José D. P. Rolim and Salil P. Vadhan, editors, *RANDOM*, volume 2483 of *LNCS*, pages 1–10. Springer, 2002. doi:10.1007/3-540-45726-7\_1.
- 6 Ezio Bartocci, Laura Kovács, and Miroslav Stankovič. Automatic generation of moment-based invariants for Prob-solvable loops. In Yu-Fang Chen, Chih-Hong Cheng, and Javier Esparza, editors, *ATVA*, volume 11781 of *LNCS*, pages 255–276. Springer, 2019. doi:10.1007/978-3-030-31784-3\_15.
- 7 Kevin Batz, Mingshuai Chen, Sebastian Junges, Benjamin Lucien Kaminski, Joost-Pieter Katoen, and Christoph Matheja. Probabilistic program verification via inductive synthesis of inductive invariants. In Sriram Sankaranarayanan and Natasha Sharygina, editors, *TACAS*, volume 13994 of *LNCS*, pages 410–429. Springer, 2023. doi:10.1007/978-3-031-30820-8\_25.
- 8 Jarosław Błasiok. Optimal streaming and tracking distinct elements with high probability. *ACM Trans. Algorithms*, 16(1):3:1–3:28, 2020. doi:10.1145/3309193.
- 9 Azucena Garvía Bosshard, Jonathan Bootle, and Christoph Sprenger. Formal verification of the Sumcheck protocol. In *CSF*, pages 605–619. IEEE, 2024. doi:10.1109/CSF61375.2024.00014.
- 10 Sourav Chakraborty, N. V. Vinodchandran, and Kuldeep S. Meel. Distinct elements in streams: An algorithm for the (text) book. In Shiri Chechik, Gonzalo Navarro, Eva Rotenberg, and Grzegorz Herman, editors, *ESA*, volume 244 of *LIPICs*, pages 34:1–34:6. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2022. doi:10.4230/LIPICs.ESA.2022.34.
- 11 Sourav Chakraborty, N. V. Vinodchandran, and Kuldeep S. Meel. Distinct elements in streams: An algorithm for the (text) book. *CoRR*, abs/2301.10191, 2023. doi:10.48550/arXiv.2301.10191.
- 12 Markus de Medeiros, Muhammad Naveed, Tancrede Lepoint, Temesghen Kahsai, Tristan Ravitch, Stefan Zetsche, Anjali Joshi, Joseph Tassarotti, Aws Albarghouthi, and Jean-Baptiste Tristan. Verified foundations for differential privacy. *CoRR*, abs/2412.01671, 2024. doi:10.48550/arXiv.2412.01671.
- 13 Jeffrey Dean and Sanjay Ghemawat. MapReduce: a flexible data processing tool. *Commun. ACM*, 53(1):72–77, 2010. doi:10.1145/1629175.1629198.
- 14 Jerry den Hartog and Erik P. de Vink. Verifying probabilistic programs using a Hoare like logic. *Int. J. Found. Comput. Sci.*, 13(3):315–340, 2002. doi:10.1142/S012905410200114X.
- 15 Devdatt P. Dubhashi, Volker Priebe, and Desh Ranjan. Negative dependence through the FKG inequality. *BRICS Report Series*, 3(27), 1996. doi:10.7146/brics.v3i27.20008.
- 16 Devdatt P. Dubhashi and Desh Ranjan. Balls and bins: A study in negative dependence. *Random Struct. Algorithms*, 13(2):99–124, 1998. doi:10.1002/(SICI)1098-2418(199809)13:2<99::AID-RSA1>3.0.CO;2-M.
- 17 Manuel Eberl, Max W. Haslbeck, and Tobias Nipkow. Verified analysis of random binary tree structures. *J. Autom. Reason.*, 64(5):879–910, 2020. doi:10.1007/S10817-020-09545-0.
- 18 Philippe Flajolet, Éric Fusy, Olivier Gandouet, and Frédéric Meunier. HyperLogLog: the analysis of a near-optimal cardinality estimation algorithm. In *Conference on Analysis of Algorithms*, 2007. doi:10.46298/dmtcs.3545.
- 19 Philippe Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. *J. Comput. Syst. Sci.*, 31(2):182–209, 1985. doi:10.1016/0022-0000(85)90041-8.
- 20 C. M. Fortuin, P. W. Kasteleyn, and J. Ginibre. Correlation inequalities on some partially ordered sets. *Commun. Math. Phys.*, 22(2):89–103, 1971. doi:10.1007/BF01651330.
- 21 Michèle Giry. A categorical approach to probability theory. In B. Banaschewski, editor, *Categorical Aspects of Topology and Analysis*, volume 915 of *LNLM*, pages 68–85. Springer, 1982. doi:10.1007/BFb0092872.

- 22 Kiran Gopinathan and Ilya Sergey. Certifying certainty and uncertainty in approximate membership query structures. In Shuvendu K. Lahiri and Chao Wang, editors, *CAV*, volume 12225 of *LNCS*, pages 279–303. Springer, 2020. doi:10.1007/978-3-030-53291-8\_16.
- 23 Peter J. Haas, Jeffrey F. Naughton, S. Seshadri, and Lynne Stokes. Sampling-based estimation of the number of distinct values of an attribute. In Umeshwar Dayal, Peter M. D. Gray, and Shojiro Nishio, editors, *VLDB*, pages 311–322. Morgan Kaufmann, 1995. URL: <http://www.vldb.org/conf/1995/P311.PDF>.
- 24 Philipp G. Haselwarter, Kwing Hei Li, Alejandro Aguirre, Simon Oddershede Gregersen, Joseph Tassarotti, and Lars Birkedal. Approximate relational reasoning for higher-order probabilistic programs. *Proc. ACM Program. Lang.*, 9(POPL):1196–1226, 2025. doi:10.1145/3704877.
- 25 Maximilian P. L. Haslbeck and Tobias Nipkow. Verified analysis of list update algorithms. In Akash Lal, S. Akshay, Saket Saurabh, and Sandeep Sen, editors, *FSTTCS*, volume 65 of *LIPICs*, pages 49:1–49:15. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2016. doi:10.4230/LIPICs.FSTTCS.2016.49.
- 26 Stefan Heule, Marc Nunkesser, and Alexander Hall. HyperLogLog in practice: algorithmic engineering of a state of the art cardinality estimation algorithm. In Giovanna Guerrini and Norman W. Paton, editors, *EDBT*, pages 683–692. ACM, 2013. doi:10.1145/2452376.2452456.
- 27 Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, March 1963. doi:10.2307/2282952.
- 28 Joe Hurd. Formal verification of probabilistic algorithms. Technical Report UCAM-CL-TR-566, University of Cambridge, Computer Laboratory, 2003. doi:10.48456/tr-566.
- 29 Kumar Joag-Dev and Frank Proschan. Negative association of random variables with applications. *Annals of Statistics*, 11:286–295, 1983. doi:10.1214/AOS/1176346079.
- 30 Daniel M. Kane, Jelani Nelson, and David P. Woodruff. An optimal algorithm for the distinct elements problem. In Jan Paredaens and Dirk Van Gucht, editors, *PODS*, pages 41–52. ACM, 2010. doi:10.1145/1807085.1807094.
- 31 Emin Karayel. Formalization of randomized approximation algorithms for frequency moments. In June Andronick and Leonardo de Moura, editors, *ITP*, volume 237 of *LIPICs*, pages 21:1–21:21. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2022. doi:10.4230/LIPICs.ITP.2022.21.
- 32 Emin Karayel. An embarrassingly parallel optimal-space cardinality estimation algorithm. In Nicole Megow and Adam D. Smith, editors, *APPROX/RANDOM*, volume 275 of *LIPICs*, pages 35:1–35:22. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2023. doi:10.4230/LIPICs.APPROX/RANDOM.2023.35.
- 33 Emin Karayel. Negatively associated random variables. *Archive of Formal Proofs*, January 2025. , Formal proof development. URL: [https://isa-afp.org/entries/Negative\\_Association.html](https://isa-afp.org/entries/Negative_Association.html).
- 34 Emin Karayel, Derek Khu, Kuldeep S. Meel, Yong Kiam Tan, and Seng Joe Watt. Verification of the CVM algorithm with a new recursive analysis technique. *Archive of Formal Proofs*, February 2025. , Formal proof development. URL: [https://isa-afp.org/entries/CVM\\_Distinct\\_Elements.html](https://isa-afp.org/entries/CVM_Distinct_Elements.html).
- 35 Donald E. Knuth. The CVM algorithm for estimating distinct elements in streams, 2023. Accessed: 2025-01-14. URL: <https://www-cs-faculty.stanford.edu/~knuth/papers/cvm-note.pdf>.
- 36 Andrey Kofnov, Marcel Moosbrugger, Miroslav Stankovic, Ezio Bartocci, and Efstathia Bura. Moment-based invariants for probabilistic loops with non-polynomial assignments. In Erika Ábrahám and Marco Paolieri, editors, *QEST*, volume 13479 of *LNCS*, pages 3–25. Springer, 2022. doi:10.1007/978-3-031-16336-4\_1.
- 37 Andreas Lochbihler. Probabilistic functions and cryptographic oracles in higher order logic. In Peter Thiemann, editor, *ESOP*, volume 9632 of *LNCS*, pages 503–531. Springer, 2016. doi:10.1007/978-3-662-49498-1\_20.

- 38 Annabelle McIver and Carroll Morgan. *Probabilistic loops: Invariants and variants*, pages 37–78. Springer, New York, NY, 2005. doi:10.1007/0-387-27006-X\_2.
- 39 Michael Mitzenmacher and Eli Upfal. *Probability and Computing: Randomization and Probabilistic Techniques in Algorithms and Data Analysis*. Cambridge University Press, USA, 2nd edition, 2017.
- 40 Rajeev Motwani and Prabhakar Raghavan. *Randomized Algorithms*. Cambridge University Press, 1995. doi:10.1017/CB09780511814075.
- 41 Steve Nadis. Computer scientists invent an efficient new way to count, 2024. Accessed: 2025-01-14. URL: <https://www.quantamagazine.org/computer-scientists-invent-an-efficient-new-way-to-count-20240516>.
- 42 Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. *Isabelle/HOL - A Proof Assistant for Higher-Order Logic*, volume 2283 of *LNCS*. Springer, 2002. doi:10.1007/3-540-45949-9.
- 43 Seth Pettie and Dingyu Wang. Information theoretic limits of cardinality estimation: Fisher meets Shannon. In Samir Khuller and Virginia Vassilevska Williams, editors, *STOC*, pages 556–569. ACM, 2021. doi:10.1145/3406325.3451032.
- 44 Robert H. Morris Sr. Counting large numbers of events in small registers. *Commun. ACM*, 21(10):840–842, 1978. doi:10.1145/359619.359627.
- 45 Daniel Stüwe and Manuel Eberl. Probabilistic primality testing. *Archive of Formal Proofs*, 2019. , Formal proof development. URL: [https://isa-afp.org/entries/Probabilistic\\_Prime\\_Tests.html](https://isa-afp.org/entries/Probabilistic_Prime_Tests.html).
- 46 Yong Kiam Tan, Jiong Yang, Mate Soos, Magnus O. Myreen, and Kuldeep S. Meel. Formally certified approximate model counting. In Arie Gurfinkel and Vijay Ganesh, editors, *CAV*, volume 14681 of *LNCS*, pages 153–177. Springer, 2024. doi:10.1007/978-3-031-65627-9\_8.
- 47 Joseph Tassarotti, Koundinya Vajjha, Anindya Banerjee, and Jean-Baptiste Tristan. A formal proof of PAC learnability for decision stumps. In Catalin Hritcu and Andrei Popescu, editors, *CPP*, pages 5–17. ACM, 2021. doi:10.1145/3437992.3439917.
- 48 Seng Joe Watt, Derek Khu, Emin Karayel, Kuldeep S. Meel, and Yong Kiam Tan. Verification of the CVM algorithm with the transformation-based proof. URL: [https://archive.softwareheritage.org/swh:1:dir:7df92f0347e7bd150efef42e3edbbde0037498cc;origin=https://github.com/joewatt95/CVM;visit=swh:1:snp:90c68a9518ce6768c54e4596003838d7d6579cc3;anchor=swh:1:rev:6bba527562c6911c664586df83e2fdeeadbd869a;path=/isabelle/CVM\\_Transforms/](https://archive.softwareheritage.org/swh:1:dir:7df92f0347e7bd150efef42e3edbbde0037498cc;origin=https://github.com/joewatt95/CVM;visit=swh:1:snp:90c68a9518ce6768c54e4596003838d7d6579cc3;anchor=swh:1:rev:6bba527562c6911c664586df83e2fdeeadbd869a;path=/isabelle/CVM_Transforms/).