

# Distributed Distinct Elements

Emin Karayel

March 17, 2025

## Abstract

This entry formalizes a randomized cardinality estimation data structure with asymptotically optimal space usage. It is inspired by the streaming algorithm presented by Błasiok [3] in 2018. His work closed the gap between the best-known lower bound and upper bound after a long line of research started by Flajolet and Martin [4] in 1984 and was the first to apply expander graphs (in addition to hash families) to the problem. The formalized algorithm has two improvements compared to the algorithm by Błasiok. It supports operation in parallel mode, and it relies on a simpler pseudo-random construction avoiding the use of code based extractors.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Preliminary Results</b>	<b>2</b>
<b>3</b>	<b>Blind</b>	<b>4</b>
<b>4</b>	<b>Balls and Bins</b>	<b>5</b>
<b>5</b>	<b>Tail Bounds for Expander Walks</b>	<b>11</b>
<b>6</b>	<b>Inner Algorithm</b>	<b>12</b>
<b>7</b>	<b>Accuracy without cutoff</b>	<b>18</b>
<b>8</b>	<b>Cutoff Level</b>	<b>21</b>
<b>9</b>	<b>Accuracy with cutoff</b>	<b>22</b>
<b>10</b>	<b>Outer Algorithm</b>	<b>23</b>

# 1 Introduction

The algorithm is described as functional data structures, given a seed which needs to be chosen uniformly from a initial segment of the natural numbers and globally, there are three functions:

- single - given the seed and an element from the universe computes a sketch for that singleton set
- merge - computes a sketch based on two input sketches and returns a sketch representing the union set
- estimate - computes an estimate for the cardinality of the set represented by a sketch

The main point is that a sketch requires  $\mathcal{O}(\delta^{-2} \ln(\varepsilon^{-1}) + \ln n)$  space where  $n$  is the universe size,  $\delta$  is the desired relative accuracy and  $\varepsilon$  is the desired failure probability. Note that it is easy to see that an exact solution would necessarily require  $\mathcal{O}(n)$  bits.

The algorithm is split into two parts an inner algorithm, described in Section 6, which itself is already a full cardinality estimation algorithm, however its space usage is below optimal. The outer algorithm is introduced in Section 10, which runs mutiple copies of the inner algorithm with carefully chosen inner parameters.

As mentioned in the abstract the algorithm is inspired by the solution to the streaming version of the problem by Błasiok [3] in 2020. His work builds on a long line of reasrch starting in 1985 [4, 1, 2, 7, 11, 5].

In an earlier AFP entry [9] I have formalized an earlier cardinality estimation algorithm based on the work by Bar-Yossef et al. [2] in 2002. Since then I have addressed the existence of finite fields for higher prime powers and expander graphs [8, 10]. Building on these results, the formalization of this more advanced solution presented here became possible.

The solution described here improves on the algorithms described by Błasiok in two ways (without comprising its optimal space usage). It can be used in a parallel mode of operation. Moreover the pseudo-random construction used is simpler than the solution described by Błasiok — who uses an extractor based on Parvaresh-Vardy codes [6] to sample random walks in an expander graph, which are then sub-sampled and then the walks are used to sample seeds for hash functions. In the solution presented here neither the sub-sampling step nor the extractor is needed, instead a two-stage expander construction is used, this means that the nodes of the first expander correspond to the walks in a second expander graph. The latters nodes correspond to seeds of hash functions (as in Błasiok’s solution).

The modification needed to support a parallel mode of operation is a change in the failure strategy of the solution presented in Kane et al., which is the event when the data in the sketch reequires too much space. The main issue is that in the parallel case the number of states the algorithm might reach is not bounded by the universe size and thus an estimate they make for the probability of the failure event does not transfer to the parallel case. To solve that the algorithm in this work is more conservative. Instead of failing out-right it instead increases a cutoff threshold. For which it is then possible to show an upper estimate independent of the number of reached states.

## 2 Preliminary Results

This section contains various short preliminary results used in the sections below.

**theory** *Distributed-Distinct-Elements-Preliminary*

**imports**

*Frequency-Moments.Frequency-Moments-Preliminary-Results*

*Universal-Hash-Families.Universal-Hash-Families-More-Product-PMF*

*Median-Method.Median*

*Expander-Graphs.Extra-Congruence-Method*

*Expander-Graphs.Constructive-Chernoff-Bound*

*Frequency-Moments.Landau-Ext*

*Stirling-Formula.Stirling-Formula*

**begin**

**unbundle** *intro-cong-syntax*

**lemma** *pmf-rev-mono*:

**assumes**  $\bigwedge x. x \in \text{set-pmf } p \implies x \notin Q \implies x \notin P$   
**shows**  $\text{measure } p \ P \leq \text{measure } p \ Q$   
*<proof>*

**lemma** *pmf-exp-mono*:

**fixes**  $f \ g :: 'a \Rightarrow \text{real}$   
**assumes**  $\text{integrable } (\text{measure-pmf } p) \ f \ \text{integrable } (\text{measure-pmf } p) \ g$   
**assumes**  $\bigwedge x. x \in \text{set-pmf } p \implies f \ x \leq g \ x$   
**shows**  $\text{integral}^L (\text{measure-pmf } p) \ f \leq \text{integral}^L (\text{measure-pmf } p) \ g$   
*<proof>*

**lemma** *pmf-markov*:

**assumes**  $\text{integrable } (\text{measure-pmf } p) \ f \ c > 0$   
**assumes**  $\bigwedge x. x \in \text{set-pmf } p \implies f \ x \geq 0$   
**shows**  $\text{measure } p \ \{\omega. f \ \omega \geq c\} \leq (\int \omega. f \ \omega \ \partial p) / c$  (**is** ?L ≤ ?R)  
*<proof>*

**lemma** *pair-pmf-prob-left*:

$\text{measure-pmf.prob } (\text{pair-pmf } A \ B) \ \{\omega. P \ (fst \ \omega)\} = \text{measure-pmf.prob } A \ \{\omega. P \ \omega\}$  (**is** ?L = ?R)  
*<proof>*

**lemma** *pmf-exp-of-fin-function*:

**assumes**  $\text{finite } A \ g \ ' \ \text{set-pmf } p \subseteq A$   
**shows**  $(\int \omega. f \ (g \ \omega) \ \partial p) = (\sum y \in A. f \ y * \text{measure } p \ \{\omega. g \ \omega = y\})$   
(**is** ?L = ?R)  
*<proof>*

Cardinality rules for distinct/ordered pairs of a set without the finiteness constraint - to use in simplification:

**lemma** *card-distinct-pairs*:

$\text{card } \{x \in B \times B. fst \ x \neq snd \ x\} = \text{card } B^{\wedge} 2 - \text{card } B$  (**is** card ?L = ?R)  
*<proof>*  
**include** *intro-cong-syntax*  
*<proof>*

**lemma** *card-ordered-pairs'*:

**fixes**  $M :: ('a :: \text{linorder}) \ \text{set}$   
**shows**  $\text{card } \{(x,y) \in M \times M. x < y\} = \text{card } M * (\text{card } M - 1) / 2$   
*<proof>*

The following are versions of the mean value theorem, where the interval endpoints may be reversed.

**lemma** *MVT-symmetric*:

**assumes**  $\bigwedge x. \llbracket \min \ a \ b \leq x; x \leq \max \ a \ b \rrbracket \implies \text{DERIV } f \ x \ :> f' \ x$   
**shows**  $\exists z :: \text{real}. \min \ a \ b \leq z \wedge z \leq \max \ a \ b \wedge (f \ b - f \ a = (b - a) * f' \ z)$   
*<proof>*

**lemma** *MVT-interval*:

**fixes**  $I :: \text{real set}$   
**assumes**  $\text{interval } I \ a \in I \ b \in I$   
**assumes**  $\bigwedge x. x \in I \implies \text{DERIV } f \ x \ :> f' \ x$   
**shows**  $\exists z. z \in I \wedge (f \ b - f \ a = (b - a) * f' \ z)$   
*<proof>*

Ln is monotone on the positive numbers and thus commutes with min and max:

**lemma** *ln-min-swap*:

$x > (0::real) \implies (y > 0) \implies \ln (\min x y) = \min (\ln x) (\ln y)$   
*<proof>*

**lemma** *ln-max-swap*:

$x > (0::real) \implies (y > 0) \implies \ln (\max x y) = \max (\ln x) (\ln y)$   
*<proof>*

Loose lower bounds for the factorial function:

**lemma** *fact-lower-bound*:

$\sqrt{2\pi n} (n/\exp(1))^n \leq \text{fact } n$  (**is** ?L ≤ ?R)  
*<proof>*

**lemma** *fact-lower-bound-1*:

**assumes**  $n > 0$   
**shows**  $(n/\exp 1)^n \leq \text{fact } n$  (**is** ?L ≤ ?R)  
*<proof>*

Rules to handle O-notation with multiple variables, where some filters may be towards zero:

**lemma** *real-inv-at-right-0-inf*:

$\forall_F x$  in *at-right*  $(0::real)$ .  $c \leq 1 / x$   
*<proof>*

**lemma** *bigO-prod-1*:

**assumes**  $(\lambda x. f x) \in O[F](\lambda x. g x)$   $G \neq \text{bot}$   
**shows**  $(\lambda x. f (fst x)) \in O[F \times_F G](\lambda x. g (fst x))$   
*<proof>*

**lemma** *bigO-prod-2*:

**assumes**  $(\lambda x. f x) \in O[G](\lambda x. g x)$   $F \neq \text{bot}$   
**shows**  $(\lambda x. f (snd x)) \in O[F \times_F G](\lambda x. g (snd x))$   
*<proof>*

**lemma** *eventually-inv*:

**fixes**  $P :: real \Rightarrow bool$   
**assumes** *eventually*  $(\lambda x. P (1/x))$  *at-top*  
**shows** *eventually*  $(\lambda x. P x)$  (*at-right* 0)  
*<proof>*

**lemma** *bigO-inv*:

**fixes**  $f g :: real \Rightarrow real$   
**assumes**  $(\lambda x. f (1/x)) \in O(\lambda x. g (1/x))$   
**shows**  $f \in O[\text{at-right } 0](g)$   
*<proof>*

**unbundle** *no intro-cong-syntax*

### 3 Blind

Blind section added to preserve section numbers

**end**

## 4 Balls and Bins

The balls and bins model describes the probability space of throwing  $r$  balls into  $b$  bins. This section derives the expected number of bins hit by at least one ball, as well as the variance in the case that each ball is thrown independently. Further, using an approximation argument it is then possible to derive bounds for the same measures in the case when the balls are being thrown only  $k$ -wise independently. The proofs follow the reasoning described in [7, §A.1] but improve on the constants, as well as constraints.

**theory** *Distributed-Distinct-Elements-Balls-and-Bins*

**imports**

*Distributed-Distinct-Elements-Preliminary*

*Discrete-Summation.Factorials*

*HOL-Combinatorics.Stirling*

*HOL-Computational-Algebra.Polynomial*

*HOL-Decision-Procs.Approximation*

**begin**

**hide-fact** *Henstock-Kurzweil-Integration.integral-sum*

**hide-fact** *Henstock-Kurzweil-Integration.integral-mult-right*

**hide-fact** *Henstock-Kurzweil-Integration.integral-nonneg*

**hide-fact** *Henstock-Kurzweil-Integration.integral-cong*

**unbundle** *intro-cong-syntax*

**lemma** *sum-power-distrib:*

**fixes**  $f :: 'a \Rightarrow \text{real}$

**assumes** *finite R*

**shows**  $(\sum_{i \in R}. f\ i) \wedge^s = (\sum xs \mid \text{set } xs \subseteq R \wedge \text{length } xs = s. (\prod x \leftarrow xs. f\ x))$

*<proof>*

**lemma** *sum-telescope-eq:*

**fixes**  $f :: \text{nat} \Rightarrow 'a :: \{\text{comm-ring-1}\}$

**shows**  $(\sum_{k \in \{\text{Suc } m..n\}}. f\ k - f\ (k - 1)) = \text{of-bool}(m \leq n) * (f\ n - f\ m)$

*<proof>*

An improved version of *diff-power-eq-sum*.

**lemma** *power-diff-sum:*

**fixes**  $a\ b :: 'a :: \{\text{comm-ring-1}, \text{power}\}$

**shows**  $a \wedge^k - b \wedge^k = (a - b) * (\sum_{i = 0..<k}. a \wedge^i * b \wedge^{(k - 1 - i)})$

*<proof>*

**lemma** *power-diff-est:*

**assumes**  $(a :: \text{real}) \geq b$

**assumes**  $b \geq 0$

**shows**  $a \wedge^k - b \wedge^k \leq (a - b) * k * a \wedge^{(k-1)}$

*<proof>*

**lemma** *power-diff-est-2:*

**assumes**  $(a :: \text{real}) \geq b$

**assumes**  $b \geq 0$

**shows**  $a \wedge^k - b \wedge^k \geq (a - b) * k * b \wedge^{(k-1)}$

*<proof>*

**lemma** *of-bool-prod:*

**assumes** *finite R*

**shows**  $(\prod_{j \in R}. \text{of-bool}(f\ j)) = (\text{of-bool}(\forall j \in R. f\ j) :: \text{real})$

*<proof>*

Additional results about falling factorials:

**lemma** *ffact-nonneg*:

**fixes**  $x :: \text{real}$   
**assumes**  $k - 1 \leq x$   
**shows**  $\text{ffact } k \ x \geq 0$   
 $\langle \text{proof} \rangle$

**lemma** *ffact-pos*:

**fixes**  $x :: \text{real}$   
**assumes**  $k - 1 < x$   
**shows**  $\text{ffact } k \ x > 0$   
 $\langle \text{proof} \rangle$

**lemma** *ffact-mono*:

**fixes**  $x \ y :: \text{real}$   
**assumes**  $k - 1 \leq x \leq y$   
**shows**  $\text{ffact } k \ x \leq \text{ffact } k \ y$   
 $\langle \text{proof} \rangle$

**lemma** *ffact-of-nat-nonneg*:

**fixes**  $x :: 'a :: \{\text{comm-ring-1}, \text{linordered-nonzero-semiring}\}$   
**assumes**  $x \in \mathbb{N}$   
**shows**  $\text{ffact } k \ x \geq 0$   
 $\langle \text{proof} \rangle$

**lemma** *ffact-suc-diff*:

**fixes**  $x :: ('a :: \text{comm-ring-1})$   
**shows**  $\text{ffact } k \ x - \text{ffact } k \ (x - 1) = \text{of-nat } k * \text{ffact } (k - 1) \ (x - 1)$  (is ?L = ?R)  
 $\langle \text{proof} \rangle$

**lemma** *ffact-bound*:

$\text{ffact } k \ (n :: \text{nat}) \leq n^k$   
 $\langle \text{proof} \rangle$

**lemma** *fact-moment-binomial*:

**fixes**  $n :: \text{nat}$  **and**  $\alpha :: \text{real}$   
**assumes**  $\alpha \in \{0..1\}$   
**defines**  $p \equiv \text{binomial-pmf } n \ \alpha$   
**shows**  $(\int \omega. \text{ffact } s \ (\text{real } \omega) \ \partial p) = \text{ffact } s \ (\text{real } n) * \alpha^s$  (is ?L = ?R)  
 $\langle \text{proof} \rangle$

The following describes polynomials of a given maximal degree as a subset of the functions, similar to the subsets  $\mathbb{Z}$  or  $\mathbb{Q}$  as subsets of larger number classes.

**definition** *Polynomials* ( $\langle \mathbb{P} \rangle$ )

**where**  $\text{Polynomials } k = \{f. \exists p. f = \text{poly } p \wedge \text{degree } p \leq k\}$

**lemma** *Polynomials-mono*:

**assumes**  $s \leq t$   
**shows**  $\mathbb{P} \ s \subseteq \mathbb{P} \ t$   
 $\langle \text{proof} \rangle$

**lemma** *Polynomials-addI*:

**assumes**  $f \in \mathbb{P} \ k \ g \in \mathbb{P} \ k$   
**shows**  $(\lambda \omega. f \ \omega + g \ \omega) \in \mathbb{P} \ k$   
 $\langle \text{proof} \rangle$

**lemma** *Polynomials-diffI*:

**fixes**  $f \ g :: 'a :: \text{comm-ring} \Rightarrow 'a$

**assumes**  $f \in \mathbb{P} k \ g \in \mathbb{P} k$   
**shows**  $(\lambda x. f x - g x) \in \mathbb{P} k$   
 $\langle proof \rangle$

**lemma** *Polynomials-idI*:  
 $(\lambda x. x) \in (\mathbb{P} 1 :: ('a :: comm-ring-1 \Rightarrow 'a) set)$   
 $\langle proof \rangle$

**lemma** *Polynomials-constI*:  
 $(\lambda x. c) \in \mathbb{P} k$   
 $\langle proof \rangle$

**lemma** *Polynomials-multI*:  
**fixes**  $f g :: 'a :: \{comm-ring\} \Rightarrow 'a$   
**assumes**  $f \in \mathbb{P} s \ g \in \mathbb{P} t$   
**shows**  $(\lambda x. f x * g x) \in \mathbb{P} (s+t)$   
 $\langle proof \rangle$

**lemma** *Polynomials-composeI*:  
**fixes**  $f g :: 'a :: \{comm-semiring-0, semiring-no-zero-divisors\} \Rightarrow 'a$   
**assumes**  $f \in \mathbb{P} s \ g \in \mathbb{P} t$   
**shows**  $(\lambda x. f (g x)) \in \mathbb{P} (s*t)$   
 $\langle proof \rangle$

**lemma** *Polynomials-const-left-multI*:  
**fixes**  $c :: 'a :: \{comm-ring\}$   
**assumes**  $f \in \mathbb{P} k$   
**shows**  $(\lambda x. c * f x) \in \mathbb{P} k$   
 $\langle proof \rangle$

**lemma** *Polynomials-const-right-multI*:  
**fixes**  $c :: 'a :: \{comm-ring\}$   
**assumes**  $f \in \mathbb{P} k$   
**shows**  $(\lambda x. f x * c) \in \mathbb{P} k$   
 $\langle proof \rangle$

**lemma** *Polynomials-const-divI*:  
**fixes**  $c :: 'a :: \{field\}$   
**assumes**  $f \in \mathbb{P} k$   
**shows**  $(\lambda x. f x / c) \in \mathbb{P} k$   
 $\langle proof \rangle$

**lemma** *Polynomials-ffact*:  $(\lambda x. ffact s (x - y)) \in (\mathbb{P} s :: ('a :: comm-ring-1 \Rightarrow 'a) set)$   
 $\langle proof \rangle$

**lemmas** *Polynomials-intros* =  
*Polynomials-const-divI*  
*Polynomials-composeI*  
*Polynomials-const-left-multI*  
*Polynomials-const-right-multI*  
*Polynomials-multI*  
*Polynomials-addI*  
*Polynomials-diffI*  
*Polynomials-idI*  
*Polynomials-constI*  
*Polynomials-ffact*

**definition**  $C_2 :: real$  where  $C_2 = 7.5$

**definition**  $C_3 :: \text{real}$  **where**  $C_3 = 16$

A locale fixing the sets of balls and bins

**locale** *balls-and-bins-abs* =  
**fixes**  $R :: 'a \text{ set}$  **and**  $B :: 'b \text{ set}$   
**assumes** *fin-B*: *finite*  $B$  **and** *B-ne*:  $B \neq \{\}$   
**assumes** *fin-R*: *finite*  $R$   
**begin**

Independent balls and bins space:

**definition**  $\Omega$   
**where**  $\Omega = \text{prod-pmf } R \ (\lambda\cdot. \text{pmf-of-set } B)$

**lemma** *set-pmf- $\Omega$* : *set-pmf*  $\Omega = R \rightarrow_E B$   
 $\langle \text{proof} \rangle$

**lemma** *card-B-gt-0*: *card*  $B > 0$   
 $\langle \text{proof} \rangle$

**lemma** *card-B-ge-1*: *card*  $B \geq 1$   
 $\langle \text{proof} \rangle$

**definition**  $Z \ j \ \omega = \text{real} \ (\text{card} \ \{i. i \in R \wedge \omega \ i = (j::'b)\})$

**definition**  $Y \ \omega = \text{real} \ (\text{card} \ (\omega \ ' R))$

**definition**  $\mu = \text{real} \ (\text{card } B) * (1 - (1-1/\text{real} \ (\text{card } B))^{\wedge \text{card } R})$

Factorial moments for the random variable describing the number of times a bin will be hit:

**lemma** *fact-moment-balls-and-bins*:  
**assumes**  $J \subseteq B \ J \neq \{\}$   
**shows**  $(\int \omega. \text{ffact } s \ (\sum j \in J. Z \ j \ \omega) \ \partial \Omega) =$   
 $\text{ffact } s \ (\text{real} \ (\text{card } R)) * (\text{real} \ (\text{card } J) / \text{real} \ (\text{card } B))^{\wedge s}$   
**(is ?L = ?R)**  
 $\langle \text{proof} \rangle$

Expectation and variance for the number of distinct bins that are hit by at least one ball in the fully independent model. The result for the variance is improved by a factor of 4 w.r.t. the paper.

**lemma**  
**shows** *exp-balls-and-bins*: *measure-pmf.expectation*  $\Omega \ Y = \mu$  **(is ?AL = ?AR)**  
**and** *var-balls-and-bins*: *measure-pmf.variance*  $\Omega \ Y \leq \text{card } R * (\text{real} \ (\text{card } R) - 1) / \text{card } B$   
**(is ?BL  $\leq$  ?BR)**  
 $\langle \text{proof} \rangle$

**definition** *lim-balls-and-bins*  $k \ p =$  (  
*prob-space.k-wise-indep-vars* (*measure-pmf*  $p$ )  $k \ (\lambda\cdot. \text{discrete}) \ (\lambda x \ \omega. \omega \ x) \ R \wedge$   
 $(\forall x. x \in R \longrightarrow \text{map-pmf} \ (\lambda \omega. \omega \ x) \ p = \text{pmf-of-set } B)$ )

**lemma** *indep*:  
**assumes** *lim-balls-and-bins*  $k \ p$   
**shows** *prob-space.k-wise-indep-vars* (*measure-pmf*  $p$ )  $k \ (\lambda\cdot. \text{discrete}) \ (\lambda x \ \omega. \omega \ x) \ R$   
 $\langle \text{proof} \rangle$

**lemma** *ran*:  
**assumes** *lim-balls-and-bins*  $k \ p \ x \in R$   
**shows** *map-pmf*  $(\lambda \omega. \omega \ x) \ p = \text{pmf-of-set } B$   
 $\langle \text{proof} \rangle$



**lemma** *Z-integrable*:

**fixes**  $f :: \text{real} \Rightarrow \text{real}$   
**assumes** *lim-balls-and-bins*  $k\ p$   
**shows** *integrable*  $p\ (\lambda\omega. f\ (Z\ i\ \omega))$   
 $\langle \text{proof} \rangle$

**lemma** *Z-any-integrable-2*:

**fixes**  $f :: \text{real} \Rightarrow \text{real}$   
**assumes** *lim-balls-and-bins*  $k\ p$   
**shows** *integrable*  $p\ (\lambda\omega. f\ (Z\ i\ \omega + Z\ j\ \omega))$   
 $\langle \text{proof} \rangle$

**lemma** *hit-count-prod-exp*:

**assumes**  $j1 \in B\ j2 \in B\ s+t \leq k$   
**assumes** *lim-balls-and-bins*  $k\ p$   
**defines**  $L \equiv \{(xs,ys). \text{set } xs \subseteq R \wedge \text{set } ys \subseteq R \wedge$   
 $(\text{set } xs \cap \text{set } ys = \{\}) \vee j1 = j2) \wedge \text{length } xs = s \wedge \text{length } ys = t\}$   
**shows**  $(\int \omega. Z\ j1\ \omega^{\wedge s} * Z\ j2\ \omega^{\wedge t} \partial p) =$   
 $(\sum (xs,ys) \in L. (1/\text{real } (\text{card } B))^{\wedge (\text{card } (\text{set } xs \cup \text{set } ys))})$   
 $(\text{is } ?L = ?R)$   
 $\langle \text{proof} \rangle$

**lemma** *hit-count-prod-pow-eq*:

**assumes**  $i \in B\ j \in B$   
**assumes** *lim-balls-and-bins*  $k\ p$   
**assumes** *lim-balls-and-bins*  $k\ q$   
**assumes**  $s+t \leq k$   
**shows**  $(\int \omega. (Z\ i\ \omega)^{\wedge s} * (Z\ j\ \omega)^{\wedge t} \partial p) = (\int \omega. (Z\ i\ \omega)^{\wedge s} * (Z\ j\ \omega)^{\wedge t} \partial q)$   
 $\langle \text{proof} \rangle$

**lemma** *hit-count-sum-pow-eq*:

**assumes**  $i \in B\ j \in B$   
**assumes** *lim-balls-and-bins*  $k\ p$   
**assumes** *lim-balls-and-bins*  $k\ q$   
**assumes**  $s \leq k$   
**shows**  $(\int \omega. (Z\ i\ \omega + Z\ j\ \omega)^{\wedge s} \partial p) = (\int \omega. (Z\ i\ \omega + Z\ j\ \omega)^{\wedge s} \partial q)$   
 $(\text{is } ?L = ?R)$   
 $\langle \text{proof} \rangle$

**lemma** *hit-count-sum-poly-eq*:

**assumes**  $i \in B\ j \in B$   
**assumes** *lim-balls-and-bins*  $k\ p$   
**assumes** *lim-balls-and-bins*  $k\ q$   
**assumes**  $f \in \mathbb{P}\ k$   
**shows**  $(\int \omega. f\ (Z\ i\ \omega + Z\ j\ \omega) \partial p) = (\int \omega. f\ (Z\ i\ \omega + Z\ j\ \omega) \partial q)$   
 $(\text{is } ?L = ?R)$   
 $\langle \text{proof} \rangle$

**lemma** *hit-count-poly-eq*:

**assumes**  $b \in B$   
**assumes** *lim-balls-and-bins*  $k\ p$   
**assumes** *lim-balls-and-bins*  $k\ q$   
**assumes**  $f \in \mathbb{P}\ k$   
**shows**  $(\int \omega. f\ (Z\ b\ \omega) \partial p) = (\int \omega. f\ (Z\ b\ \omega) \partial q)$   $(\text{is } ?L = ?R)$   
 $\langle \text{proof} \rangle$

**lemma** *lim-balls-and-bins-from-ind-balls-and-bins*:

*lim-balls-and-bins k Ω*  
 ⟨proof⟩

**lemma** *hit-count-factorial-moments:*

**assumes**  $a:j \in B$   
**assumes**  $s \leq k$   
**assumes** *lim-balls-and-bins k p*  
**shows**  $(\int \omega. \text{ffact } s (Z j \omega) \partial p) = \text{ffact } s (\text{real } (\text{card } R)) * (1 / \text{real } (\text{card } B))^s$   
**(is ?L = ?R)**

⟨proof⟩

**lemma** *hit-count-factorial-moments-2:*

**assumes**  $a:i \in B j \in B$   
**assumes**  $i \neq j s \leq k \text{card } R \leq \text{card } B$   
**assumes** *lim-balls-and-bins k p*  
**shows**  $(\int \omega. \text{ffact } s (Z i \omega + Z j \omega) \partial p) \leq 2^s$   
**(is ?L ≤ ?R)**

⟨proof⟩

**lemma** *balls-and-bins-approx-helper:*

**fixes**  $x :: \text{real}$   
**assumes**  $x \geq 2$   
**assumes**  $\text{real } k \geq 5*x / \ln x$   
**shows**  $k \geq 2$   
**and**  $2^{k+3} / \text{fact } k \leq (1/\exp x)^2$   
**and**  $2 / \text{fact } k \leq 1 / (\exp 1 * \exp x)$

⟨proof⟩

Bounds on the expectation and variance in the k-wise independent case. Here the independence assumption is improved by a factor of two compared to the result in the paper.

**lemma**

**assumes**  $\text{card } R \leq \text{card } B$   
**assumes**  $\bigwedge c. \text{lim-balls-and-bins } (k+1) (p c)$   
**assumes**  $\varepsilon \in \{0 < .. 1/\exp(2)\}$   
**assumes**  $k \geq 5 * \ln (\text{card } B / \varepsilon) / \ln (\ln (\text{card } B / \varepsilon))$

**shows**

*exp-approx:*  $|\text{measure-pmf.expectation } (p \text{ True}) Y - \text{measure-pmf.expectation } (p \text{ False}) Y| \leq \varepsilon * \text{real } (\text{card } R)$  **(is ?A) and**  
*var-approx:*  $|\text{measure-pmf.variance } (p \text{ True}) Y - \text{measure-pmf.variance } (p \text{ False}) Y| \leq \varepsilon^2$   
**(is ?B)**

⟨proof⟩

**lemma**

**assumes**  $\text{card } R \leq \text{card } B$   
**assumes** *lim-balls-and-bins (k+1) p*  
**assumes**  $k \geq 7.5 * (\ln (\text{card } B) + 2)$   
**shows** *exp-approx-2:*  $|\text{measure-pmf.expectation } p Y - \mu| \leq \text{card } R / \text{sqrt } (\text{card } B)$   
**(is ?AL ≤ ?AR)**  
**and** *var-approx-2:*  $\text{measure-pmf.variance } p Y \leq \text{real } (\text{card } R)^2 / \text{card } B$   
**(is ?BL ≤ ?BR)**

⟨proof⟩

**lemma** *deviation-bound:*

**assumes**  $\text{card } R \leq \text{card } B$   
**assumes** *lim-balls-and-bins k p*  
**assumes**  $\text{real } k \geq C_2 * \ln (\text{real } (\text{card } B)) + C_3$   
**shows**  $\text{measure } p \{\omega. |Y \omega - \mu| > 9 * \text{real } (\text{card } R) / \text{sqrt } (\text{real } (\text{card } B))\} \leq 1 / 2^6$   
**(is ?L ≤ ?R)**

⟨proof⟩

end

unbundle *no intro-cong-syntax*

end

## 5 Tail Bounds for Expander Walks

**theory** *Distributed-Distinct-Elements-Tail-Bounds*

**imports**

*Distributed-Distinct-Elements-Preliminary*

*Expander-Graphs.Pseudorandom-Objects-Expander-Walks*

*HOL-Decision-Procs.Approximation*

**begin**

This section introduces tail estimates for random walks in expander graphs, specific to the verification of this algorithm (in particular to two-stage expander graph sampling and obtained tail bounds for subgaussian random variables). They follow from the more fundamental results *regular-graph.kl-chernoff-property* and *regular-graph.uniform-property* which are verified in the AFP entry for expander graphs [10].

**hide-fact** *Henstock-Kurzweil-Integration.integral-sum*

unbundle *intro-cong-syntax*

**lemma** *x-ln-x-min:*

**assumes**  $x \geq (0::real)$

**shows**  $x * \ln x \geq -exp(-1)$

⟨proof⟩

**theorem** (in *regular-graph*) *walk-tail-bound:*

**assumes**  $l > 0$

**assumes**  $S \subseteq \text{verts } G$

**defines**  $\mu \equiv \text{real } (\text{card } S) / \text{card } (\text{verts } G)$

**assumes**  $\gamma < 1 \ \mu + \Lambda_a \leq \gamma$

**shows**  $\text{measure } (\text{pmf-of-multiset } (\text{walks } G \ l)) \{w. \text{real } (\text{card } \{i \in \{..<l\}. w ! i \in S\}) \geq \gamma * l\} \\ \leq \text{exp } (- \text{real } l * (\gamma * \ln (1/(\mu + \Lambda_a)) - 2 * \text{exp}(-1))) \text{ (is ?L \le ?R)}$

⟨proof⟩

**theorem** (in *regular-graph*) *walk-tail-bound-2:*

**assumes**  $l > 0 \ \Lambda_a \leq \Lambda \ \Lambda > 0$

**assumes**  $S \subseteq \text{verts } G$

**defines**  $\mu \equiv \text{real } (\text{card } S) / \text{card } (\text{verts } G)$

**assumes**  $\gamma < 1 \ \mu + \Lambda \leq \gamma$

**shows**  $\text{measure } (\text{pmf-of-multiset } (\text{walks } G \ l)) \{w. \text{real } (\text{card } \{i \in \{..<l\}. w ! i \in S\}) \geq \gamma * l\} \\ \leq \text{exp } (- \text{real } l * (\gamma * \ln (1/(\mu + \Lambda)) - 2 * \text{exp}(-1))) \text{ (is ?L \le ?R)}$

⟨proof⟩

**lemma** *disjI-safe:*  $(\neg x \implies y) \implies x \vee y$  ⟨proof⟩

**lemma** *walk-tail-bound:*

**fixes**  $T$

**assumes**  $l > 0 \ \Lambda > 0$

**assumes**  $\text{measure } (\text{sample-pro } S) \{w. T \ w\} \leq \mu$

**assumes**  $\gamma \leq 1 \ \mu + \Lambda \leq \gamma \ \mu \leq 1$

**shows**  $\text{measure } (\text{sample-pro } (\mathcal{E} \ l \ \Lambda \ S)) \{w. \text{real } (\text{card } \{i \in \{..<l\}. T \ (w \ i)\}) \geq \gamma * l\}$

$\leq \exp(-\text{real } l * (\gamma * \ln(1/(\mu + \Lambda)) - 2 * \exp(-1)))$  (is ?L ≤ ?R)  
 <proof>

**definition**  $C_1 :: \text{real where } C_1 = \exp 2 + \exp 3 + (\exp 1 - 1)$

**lemma** *deviation-bound:*

**fixes**  $f :: 'a \Rightarrow \text{real}$   
**assumes**  $l > 0$   
**assumes**  $\Lambda \in \{0 < .. \exp(-\text{real } l * \ln(\text{real } l)^3)\}$   
**assumes**  $\bigwedge x. x \geq 20 \implies \text{measure}(\text{sample-pro } S) \{v. f v \geq x\} \leq \exp(-x * \ln x^3)$   
**shows**  $\text{measure}(\text{sample-pro } (\mathcal{E} l \Lambda S)) \{\omega. (\sum_{i < l} f(\omega i)) \geq C_1 * l\} \leq \exp(-\text{real } l)$  (is ?L ≤ ?R)  
 <proof>

**unbundle** *no intro-cong-syntax*

**end**

## 6 Inner Algorithm

This section introduces the inner algorithm (as mentioned it is already a solution to the cardinality estimation with the caveat that, if  $\varepsilon$  is too small it requires too much space. The outer algorithm in Section 10 resolved this problem.

The algorithm makes use of the balls and bins model, more precisely, the fact that the number of hit bins can be used to estimate the number of balls thrown (even if there are collisions). I.e. it assigns each universe element to a bin using a  $k$ -wise independent hash function. Then it counts the number of bins hit.

This strategy however would only work if the number of balls is roughly equal to the number of bins, to remedy that the algorithm performs an adaptive sub-sampling strategy. This works by assigning each universe element a level (using a second hash function) with a geometric distribution. The algorithm then selects a level that is appropriate based on a rough estimate obtained using the maximum level in the bins.

To save space the algorithm drops information about small levels, whenever the space usage would be too high otherwise. This level will be called the cutoff-level. This is okay as long as the cutoff level is not larger than the sub-sampling threshold. A lot of the complexity in the proof is devoted to verifying that the cutoff-level will not cross it, it works by defining a third value  $s_M$  that is both an upper bound for the cutoff level and a lower bound for the subsampling threshold simultaneously with high probability.

**theory** *Distributed-Distinct-Elements-Inner-Algorithm*

**imports**

*Universal-Hash-Families.Pseudorandom-Objects-Hash-Families*  
*Distributed-Distinct-Elements-Preliminary*  
*Distributed-Distinct-Elements-Balls-and-Bins*  
*Distributed-Distinct-Elements-Tail-Bounds*  
*Prefix-Free-Code-Combinators.Prefix-Free-Code-Combinators*

**begin**

**unbundle** *intro-cong-syntax*

**hide-const** *Abstract-Rewriting.restrict*

**definition**  $C_4 :: \text{real where } C_4 = 3^2 * 2^{23}$

**definition**  $C_5 :: \text{int where } C_5 = 33$

**definition**  $C_6 :: \text{real where } C_6 = 4$

**definition**  $C_7 :: \text{nat where } C_7 = 2^5$

**locale** *inner-algorithm* =

**fixes**  $n :: \text{nat}$

**fixes**  $\delta :: \text{real}$

**fixes**  $\varepsilon :: \text{real}$

**assumes**  $n\text{-gt-0}: n > 0$

**assumes**  $\delta\text{-gt-0}: \delta > 0$  **and**  $\delta\text{-lt-1}: \delta < 1$

**assumes**  $\varepsilon\text{-gt-0}: \varepsilon > 0$  **and**  $\varepsilon\text{-lt-1}: \varepsilon < 1$

**begin**

**definition**  $b\text{-exp}$  **where**  $b\text{-exp} = \text{nat } \lceil \log 2 (C_4 / \varepsilon^2) \rceil$

**definition**  $b :: \text{nat}$  **where**  $b = 2^{b\text{-exp}}$

**definition**  $l$  **where**  $l = \text{nat } \lceil C_6 * \ln (2 / \delta) \rceil$

**definition**  $k$  **where**  $k = \text{nat } \lceil C_2 * \ln b + C_3 \rceil$

**definition**  $\Lambda :: \text{real}$  **where**  $\Lambda = \min (1/16) (\exp (-l * \ln l^3))$

**definition**  $\varrho :: \text{real} \Rightarrow \text{real}$  **where**  $\varrho x = b * (1 - (1-1/b)^{\text{powr } x})$

**definition**  $\varrho\text{-inv} :: \text{real} \Rightarrow \text{real}$  **where**  $\varrho\text{-inv } x = \ln (1-x/b) / \ln (1-1/b)$

**lemma**  $l\text{-lbound}: C_6 * \ln (2 / \delta) \leq l$

*<proof>*

**lemma**  $k\text{-min}: C_2 * \ln (\text{real } b) + C_3 \leq \text{real } k$

*<proof>*

**lemma**  $\Lambda\text{-gt-0}: \Lambda > 0$

*<proof>*

**lemma**  $\Lambda\text{-le-1}: \Lambda \leq 1$

*<proof>*

**lemma**  $l\text{-gt-0}: l > 0$

*<proof>*

**lemma**  $l\text{-ubound}: l \leq C_6 * \ln (1 / \delta) + C_6 * \ln 2 + 1$

*<proof>*

**lemma**  $b\text{-exp-ge-26}: b\text{-exp} \geq 26$

*<proof>*

**lemma**  $b\text{-min}: b \geq 2^{26}$

*<proof>*

**lemma**  $k\text{-gt-0}: k > 0$

*<proof>*

**lemma**  $b\text{-ne}: \{..<b\} \neq \{\}$

*<proof>*

**lemma**  $b\text{-lower-bound}: C_4 / \varepsilon^2 \leq \text{real } b$

*<proof>*

**definition**  $n\text{-exp}$  **where**  $n\text{-exp} = \max (\text{nat } \lceil \log 2 n \rceil) 1$

**lemma**  $n\text{-exp-gt-0}: n\text{-exp} > 0$

*<proof>*

**abbreviation**  $\Psi_1$  **where**  $\Psi_1 \equiv \mathcal{H} \ 2 \ n \ (\mathcal{G} \ n\text{-exp})$

**abbreviation**  $\Psi_2$  **where**  $\Psi_2 \equiv \mathcal{H} \ 2 \ n \ (\mathcal{N} \ (C_7 * b^2))$

**abbreviation**  $\Psi_3$  **where**  $\Psi_3 \equiv \mathcal{H} \ k \ (C_7 * b^2) \ (\mathcal{N} \ b)$

**definition**  $\Psi$  **where**  $\Psi = \Psi_1 \times_P \Psi_2 \times_P \Psi_3$

**abbreviation**  $\Omega$  **where**  $\Omega \equiv \mathcal{E} \ l \ \Lambda \ \Psi$

**type-synonym**  $state = (nat \Rightarrow nat \Rightarrow int) \times (nat)$

**fun**  $is\text{-}too\text{-}large :: (nat \Rightarrow nat \Rightarrow int) \Rightarrow bool$  **where**  
 $is\text{-}too\text{-}large \ B = ((\sum (i,j) \in \{..\<l\} \times \{..\<b\}. \lfloor \log 2 (\max (B \ i \ j) \ (-1) + 2) \rfloor) > C_5 * b * l)$

**fun**  $compress\text{-}step :: state \Rightarrow state$  **where**  
 $compress\text{-}step \ (B,q) = (\lambda \ i \ j. \max (B \ i \ j - 1) \ (-1), \ q+1)$

**function**  $compress :: state \Rightarrow state$  **where**

$compress \ (B,q) = ($   
   $if \ is\text{-}too\text{-}large \ B$   
     $then \ (compress \ (compress\text{-}step \ (B,q)))$   
     $else \ (B,q)$   
 $\langle proof \rangle$

**fun**  $compress\text{-}termination :: state \Rightarrow nat$  **where**  
 $compress\text{-}termination \ (B,q) = (\sum (i,j) \in \{..\<l\} \times \{..\<b\}. \ nat \ (B \ i \ j + 1))$

**lemma**  $compress\text{-}termination$ :

**assumes**  $is\text{-}too\text{-}large \ B$

**shows**  $compress\text{-}termination \ (compress\text{-}step \ (B,q)) < compress\text{-}termination \ (B,q)$

$\langle proof \rangle$

**termination**  $compress$

$\langle proof \rangle$

**fun**  $merge1 :: state \Rightarrow state \Rightarrow state$  **where**

$merge1 \ (B1,q1) \ (B2, \ q2) = ($   
   $let \ q = \max \ q1 \ q2 \ in \ (\lambda \ i \ j. \max (B1 \ i \ j + q1 - q) \ (B2 \ i \ j + q2 - q), \ q)$

**fun**  $merge :: state \Rightarrow state \Rightarrow state$  **where**

$merge \ x \ y = compress \ (merge1 \ x \ y)$

**type-synonym**  $seed = nat \Rightarrow (nat \Rightarrow nat) \times (nat \Rightarrow nat) \times (nat \Rightarrow nat)$

**fun**  $single1 :: seed \Rightarrow nat \Rightarrow state$  **where**

$single1 \ \omega \ x = (\lambda \ i \ j.$   
   $let \ (f,g,h) = \omega \ i \ in \ ($   
     $if \ h \ (g \ x) = j \ \wedge \ i < l \ then \ int \ (f \ x) \ else \ (-1), \ 0)$

**fun**  $single :: seed \Rightarrow nat \Rightarrow state$  **where**

$single \ \omega \ x = compress \ (single1 \ \omega \ x)$

**fun**  $estimate1 :: state \Rightarrow nat \Rightarrow real$  **where**

$estimate1 \ (B,q) \ i = ($   
   $let \ s = \max \ 0 \ (\Max \ ((B \ i) \ ' \ \{..\<b\}) + q - \lfloor \log 2 \ b \rfloor + 9);$   
   $p = \text{card} \ \{ j. \ j \in \{..\<b\} \ \wedge \ B \ i \ j + q \geq s \} \ in$   
   $2 \ \text{powl} \ s * \ln \ (1-p/b) / \ln(1-1/b))$

**fun**  $estimate :: state \Rightarrow real$  **where**

$estimate \ x = \text{median} \ l \ (estimate1 \ x)$

## 6.1 History Independence

**fun**  $\tau_0 :: ((nat \Rightarrow nat) \times (nat \Rightarrow nat) \times (nat \Rightarrow nat)) \Rightarrow nat \ set \Rightarrow nat \Rightarrow int$   
**where**  $\tau_0 (f,g,h) A j = Max (\{ int (f a) \mid a . a \in A \wedge h (g a) = j \} \cup \{-1\})$

**definition**  $\tau_1 :: ((nat \Rightarrow nat) \times (nat \Rightarrow nat) \times (nat \Rightarrow nat)) \Rightarrow nat \ set \Rightarrow nat \Rightarrow nat \Rightarrow int$   
**where**  $\tau_1 \psi A q j = max (\tau_0 \psi A j - q) (-1)$

**definition**  $\tau_2 :: seed \Rightarrow nat \ set \Rightarrow nat \Rightarrow nat \Rightarrow nat \Rightarrow int$   
**where**  $\tau_2 \omega A q i j = (if i < l then \tau_1 (\omega i) A q j else (-1))$

**definition**  $\tau_3 :: seed \Rightarrow nat \ set \Rightarrow nat \Rightarrow state$   
**where**  $\tau_3 \omega A q = (\tau_2 \omega A q, q)$

**definition**  $q :: seed \Rightarrow nat \ set \Rightarrow nat$   
**where**  $q \omega A = (LEAST q . \neg(is-too-large (\tau_2 \omega A q)))$

**definition**  $\tau :: seed \Rightarrow nat \ set \Rightarrow state$   
**where**  $\tau \omega A = \tau_3 \omega A (q \omega A)$

**lemma**  $\tau_2$ -step:  $\tau_2 \omega A (x+y) = (\lambda i j. max (\tau_2 \omega A x i j - y) (-1))$   
 $\langle proof \rangle$

**lemma**  $\tau_3$ -step: compress-step  $(\tau_3 \omega A x) = \tau_3 \omega A (x+1)$   
 $\langle proof \rangle$

**lemma**  $\Psi_1$ : is-prime-power (pro-size ( $\mathcal{G}$  n-exp))  
 $\langle proof \rangle$

**lemma**  $\Psi_2$ : is-prime-power (pro-size ( $\mathcal{N} (C_7 * b^2)$ ))  
 $\langle proof \rangle$

**lemma**  $\Psi_3$ : is-prime-power (pro-size ( $\mathcal{N} b$ ))  
 $\langle proof \rangle$

**lemma** sample-pro- $\Psi$ :  
 $sample-pro \Psi = pair-pmf (sample-pro \Psi_1) (pair-pmf (sample-pro \Psi_2) (sample-pro \Psi_3))$   
 $\langle proof \rangle$

**lemma** sample-set- $\Psi$ : pro-set  $\Psi = pro-set \Psi_1 \times pro-set \Psi_2 \times pro-set \Psi_3$   
 $\langle proof \rangle$

**lemma** f-range:  
**assumes**  $(f,g,h) \in pro-set \Psi$   
**shows**  $f x \leq n-exp$   
 $\langle proof \rangle$

**lemma** g-range-1:  
**assumes**  $g \in pro-set \Psi_2$   
**shows**  $g x < C_7 * b^2$   
 $\langle proof \rangle$

**lemma** h-range-1:  
**assumes**  $h \in pro-set \Psi_3$   
**shows**  $h x < b$   
 $\langle proof \rangle$

**lemma** g-range:

**assumes**  $(f,g,h) \in \text{pro-set } \Psi$   
**shows**  $g\ x < C_7 * b^{\wedge} 2$   
 $\langle \text{proof} \rangle$

**lemma** *h-range*:

**assumes**  $(f,g,h) \in \text{pro-set } \Psi$   
**shows**  $h\ x < b$   
 $\langle \text{proof} \rangle$

**lemma** *fin-f*:

**assumes**  $(f,g,h) \in \text{pro-set } \Psi$   
**shows**  $\text{finite } \{ \text{int } (f\ a) \mid a. P\ a \}$  (**is finite** ?M)  
 $\langle \text{proof} \rangle$

**lemma** *Max-int-range*:  $x \leq (y::\text{int}) \implies \text{Max } \{x..y\} = y$   
 $\langle \text{proof} \rangle$

**lemma**  $\Omega$ :  $l > 0\ \Lambda > 0$   $\langle \text{proof} \rangle$

**lemma**  $\omega$ -range:

**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\omega\ i \in \text{pro-set } \Psi$   
 $\langle \text{proof} \rangle$

**lemma** *max-q-1*:

**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\tau_2\ \omega\ A\ (\text{nat } \lceil \log 2\ n \rceil + 2)\ i\ j = (-1)$   
 $\langle \text{proof} \rangle$

**lemma** *max-q-2*:

**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\neg (\text{is-too-large } (\tau_2\ \omega\ A\ (\text{nat } \lceil \log 2\ n \rceil + 2)))$   
 $\langle \text{proof} \rangle$

**lemma** *max-s-3*:

**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $q\ \omega\ A \leq (\text{nat } \lceil \log 2\ n \rceil + 2)$   
 $\langle \text{proof} \rangle$

**lemma** *max-mono*:  $x \leq (y::'a::\text{linorder}) \implies \text{max } x\ z \leq \text{max } y\ z$   
 $\langle \text{proof} \rangle$

**lemma** *max-mono-2*:  $y \leq (z::'a::\text{linorder}) \implies \text{max } x\ y \leq \text{max } x\ z$   
 $\langle \text{proof} \rangle$

**lemma**  $\tau_0$ -mono:

**assumes**  $\psi \in \text{pro-set } \Psi$   
**assumes**  $A \subseteq B$   
**shows**  $\tau_0\ \psi\ A\ j \leq \tau_0\ \psi\ B\ j$   
 $\langle \text{proof} \rangle$

**lemma**  $\tau_2$ -mono:

**assumes**  $\omega \in \text{pro-set } \Omega$   
**assumes**  $A \subseteq B$   
**shows**  $\tau_2\ \omega\ A\ x\ i\ j \leq \tau_2\ \omega\ B\ x\ i\ j$   
 $\langle \text{proof} \rangle$

**lemma** *is-too-large-antimono*:



**assumes**  $\omega \in \text{pro-set } \Omega$   
**assumes**  $A \subseteq B$   
**assumes**  $\text{is-too-large } (\tau_2 \omega A x)$   
**shows**  $\text{is-too-large } (\tau_2 \omega B x)$   
 <proof>

**lemma** *q-compact*:  
**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\neg (\text{is-too-large } (\tau_2 \omega A (q \omega A)))$   
 <proof>

**lemma** *q-mono*:  
**assumes**  $\omega \in \text{pro-set } \Omega$   
**assumes**  $A \subseteq B$   
**shows**  $q \omega A \leq q \omega B$   
 <proof>

**lemma** *lt-s-too-large*:  $x < q \omega A \implies \text{is-too-large } (\tau_2 \omega A x)$   
 <proof>

**lemma** *compress-result-1*:  
**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\text{compress } (\tau_3 \omega A (q \omega A - i)) = \tau \omega A$   
 <proof>

**lemma** *compress-result*:  
**assumes**  $\omega \in \text{pro-set } \Omega$   
**assumes**  $x \leq q \omega A$   
**shows**  $\text{compress } (\tau_3 \omega A x) = \tau \omega A$   
 <proof>

**lemma**  $\tau_0$ -merge:  
**assumes**  $(f,g,h) \in \text{pro-set } \Psi$   
**shows**  $\tau_0 (f,g,h) (A \cup B) j = \max (\tau_0 (f,g,h) A j) (\tau_0 (f,g,h) B j)$  (**is** ?L = ?R)  
 <proof>

**lemma**  $\tau_2$ -merge:  
**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\tau_2 \omega (A \cup B) x i j = \max (\tau_2 \omega A x i j) (\tau_2 \omega B x i j)$   
 <proof>

**lemma** *merge1-result*:  
**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\text{merge1 } (\tau \omega A) (\tau \omega B) = \tau_3 \omega (A \cup B) (\max (q \omega A) (q \omega B))$   
 <proof>

**lemma** *merge-result*:  
**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\text{merge } (\tau \omega A) (\tau \omega B) = \tau \omega (A \cup B)$  (**is** ?L = ?R)  
 <proof>

**lemma** *single1-result*:  $\text{single1 } \omega x = \tau_3 \omega \{x\} 0$   
 <proof>

**lemma** *single-result*:  
**assumes**  $\omega \in \text{pro-set } \Omega$   
**shows**  $\text{single } \omega x = \tau \omega \{x\}$  (**is** ?L = ?R)  
 <proof>

## 6.2 Encoding states of the inner algorithm

**definition** *is-state-table* :: (nat × nat ⇒ int) ⇒ bool **where**  
*is-state-table* g = (range g ⊆ {-1..} ∧ g ' (-({..<l} × {..<b}))) ⊆ {-1})

Encoding for state table values:

**definition** *V<sub>e</sub>* :: int encoding  
**where** *V<sub>e</sub>* x = (if x ≥ -1 then *N<sub>e</sub>* (nat (x+1)) else None)

Encoding for state table:

**definition** *T<sub>e</sub>'* :: (nat × nat ⇒ int) encoding **where**  
*T<sub>e</sub>'* g = (  
 if *is-state-table* g  
 then (List.product [0..<l] [0..<b] →<sub>e</sub> *V<sub>e</sub>*) (restrict g ({..<l} × {..<b}))  
 else None)

**definition** *T<sub>e</sub>* :: (nat ⇒ nat ⇒ int) encoding  
**where** *T<sub>e</sub>* f = *T<sub>e</sub>'* (case-prod f)

**definition** *encode-state* :: state encoding  
**where** *encode-state* = *T<sub>e</sub>* ×<sub>e</sub> *Nb<sub>e</sub>* (nat ⌈log 2 n⌉ + 3)

**lemma** *inj-on-restrict*:  
**assumes** B ⊆ {f. f ' (- A) ⊆ {c}}  
**shows** *inj-on* (λx. restrict x A) B  
 ⟨proof⟩

**lemma** *encode-state: is-encoding encode-state*  
 ⟨proof⟩

**lemma** *state-bit-count*:  
**assumes** ω ∈ pro-set Ω  
**shows** *bit-count* (encode-state (τ ω A)) ≤ 2<sup>36</sup> \* (ln(1/δ)+1) / ε<sup>2</sup> + log 2 (log 2 n + 3)  
 (is ?L ≤ ?R)  
 ⟨proof⟩

**lemma** *random-bit-count*:  
*pro-size* Ω ≤ 2 powr (4 \* log 2 n + 48 \* (log 2 (1 / ε) + 16)<sup>2</sup> + (55 + 60 \* ln (1 / δ))<sup>3</sup>)  
 (is ?L ≤ ?R)  
 ⟨proof⟩

end

unbundle no intro-cong-syntax

end

## 7 Accuracy without cutoff

This section verifies that each of the *l* estimate have the required accuracy with high probability assuming that there was no cut-off, i.e., that *s* = 0. Section 9 will then show that this remains true as long as the cut-off is below *t f* the subsampling threshold.

**theory** *Distributed-Distinct-Elements-Accuracy-Without-Cutoff*

**imports**

*Concentration-Inequalities.Bienaymes-Identity*  
*Distributed-Distinct-Elements-Inner-Algorithm*  
*Distributed-Distinct-Elements-Balls-and-Bins*

**begin**

**no-notation** *Polynomials.var* ( $\langle X_1 \rangle$ )

**locale** *inner-algorithm-fix-A* = *inner-algorithm* +

**fixes** *A*

**assumes** *A-range*:  $A \subseteq \{..<n\}$

**assumes** *A-nonempty*:  $\{\} \neq A$

**begin**

**definition** *X* :: *nat* **where**  $X = \text{card } A$

**definition** *q-max* **where**  $q\text{-max} = \text{nat } (\lceil \log 2 X \rceil - b\text{-exp})$

**definition** *t* :: (*nat*  $\Rightarrow$  *nat*)  $\Rightarrow$  *int*  
**where**  $t f = \text{int } (\text{Max } (f \text{ ' } A)) - b\text{-exp} + 9$

**definition** *s* :: (*nat*  $\Rightarrow$  *nat*)  $\Rightarrow$  *nat*  
**where**  $s f = \text{nat } (t f)$

**definition** *R* :: (*nat*  $\Rightarrow$  *nat*)  $\Rightarrow$  *nat set*  
**where**  $R f = \{a. a \in A \wedge f a \geq s f\}$

**definition** *r* :: *nat*  $\Rightarrow$  (*nat*  $\Rightarrow$  *nat*)  $\Rightarrow$  *nat*  
**where**  $r x f = \text{card } \{a. a \in A \wedge f a \geq x\}$

**definition** *p* **where**  $p = (\lambda(f,g,h). \text{card } \{j \in \{..<b\}. \tau_1 (f,g,h) A \ 0 j \geq s f\})$

**definition** *Y* **where**  $Y = (\lambda(f,g,h). 2 \wedge s f * \varrho\text{-inv } (p (f,g,h)))$

**lemma** *fin-A*: *finite A*  
*<proof>*

**lemma** *X-le-n*:  $X \leq n$   
*<proof>*

**lemma** *X-ge-1*:  $X \geq 1$   
*<proof>*

**lemma** *of-bool-square*:  $(\text{of-bool } x)^2 = ((\text{of-bool } x)::\text{real})$   
*<proof>*

**lemma** *r-eg*:  $r x f = (\sum a \in A. (\text{of-bool } (x \leq f a) :: \text{real}))$   
*<proof>*

**lemma**  
**shows**

*r-exp*:  $(\int \omega. \text{real } (r x \omega) \partial \Psi_1) = \text{real } X * (\text{of-bool } (x \leq \max (\text{nat } \lceil \log 2 n \rceil) 1) / 2 \wedge x)$  **and**  
*r-var*:  $\text{measure-pmf.variance } \Psi_1 (\lambda \omega. \text{real } (r x \omega)) \leq (\int \omega. \text{real } (r x \omega) \partial \Psi_1)$

*<proof>*

**definition** *E<sub>1</sub>* **where**  $E_1 = (\lambda(f,g,h). 2 \text{ powr } (-t f) * X \in \{b/2 \wedge 16..b/2\})$

**lemma** *t-low*:

$\text{measure } \Psi_1 \{f. \text{of-int } (t f) < \log 2 (\text{real } X) + 1 - b\text{-exp}\} \leq 1/2 \wedge 7$  (**is ?L  $\leq$  ?R**)  
*<proof>*

**lemma** *t-high*:

*measure*  $\Psi_1 \{f. \text{of-int } (t f) > \log 2 \text{ (real } X) + 16 - b\text{-exp}\} \leq 1/2^{\wedge}7$  (**is** ?L ≤ ?R)  
<proof>

**lemma** *e-1*: *measure*  $\Psi \{\psi. \neg E_1 \psi\} \leq 1/2^{\wedge}6$   
<proof>

**definition** *E<sub>2</sub>* **where**  $E_2 = (\lambda(f,g,h). |\text{card } (R f) - X / 2^{\wedge}(s f)| \leq \varepsilon/3 * X / 2^{\wedge}(s f))$

**lemma** *e-2*: *measure*  $\Psi \{\psi. E_1 \psi \wedge \neg E_2 \psi\} \leq 1/2^{\wedge}6$  (**is** ?L ≤ ?R)  
<proof>

**definition** *E<sub>3</sub>* **where**  $E_3 = (\lambda(f,g,h). \text{inj-on } g \text{ (R f)})$

**lemma** *R-bound*:

**fixes** *f g h*

**assumes**  $E_1 (f,g,h)$

**assumes**  $E_2 (f,g,h)$

**shows**  $\text{card } (R f) \leq 2/3 * b$

<proof>

**lemma** *e-3*: *measure*  $\Psi \{\psi. E_1 \psi \wedge E_2 \psi \wedge \neg E_3 \psi\} \leq 1/2^{\wedge}6$  (**is** ?L ≤ ?R)  
<proof>

**definition** *E<sub>4</sub>* **where**  $E_4 = (\lambda(f,g,h). |p (f,g,h) - \varrho (\text{card } (R f))| \leq \varepsilon/12 * \text{card } (R f))$

**lemma** *e-4-h*:  $9 / \text{sqrt } b \leq \varepsilon / 12$   
<proof>

**lemma** *e-4*: *measure*  $\Psi \{\psi. E_1 \psi \wedge E_2 \psi \wedge E_3 \psi \wedge \neg E_4 \psi\} \leq 1/2^{\wedge}6$  (**is** ?L ≤ ?R)  
<proof>

**lemma** *ρ-inverse*:  $\varrho\text{-inv } (\varrho x) = x$   
<proof>

**lemma** *rho-mono*:

**assumes**  $x \leq y$

**shows**  $\varrho x \leq \varrho y$

<proof>

**lemma** *rho-two-thirds*:  $\varrho (2/3 * b) \leq 3/5 * b$   
<proof>

**definition** *ρ-inv'* **::**  $\text{real} \Rightarrow \text{real}$

**where**  $\varrho\text{-inv}' x = -1 / (\text{real } b * (1-x / \text{real } b) * \ln (1 - 1 / \text{real } b))$

**lemma** *ρ-inv'-bound*:

**assumes**  $x \geq 0$

**assumes**  $x \leq 59/90 * b$

**shows**  $|\varrho\text{-inv}' x| \leq 4$

<proof>

**lemma** *ρ-inv'*:

**fixes**  $x :: \text{real}$

**assumes**  $x < b$

**shows**  $\text{DERIV } \varrho\text{-inv } x :> \varrho\text{-inv}' x$

<proof>

**lemma** *accuracy-without-cutoff*:

$measure \Psi \{(f,g,h). |Y(f,g,h) - real X| > \varepsilon * X \vee s f < q-max\} \leq 1/2^4$   
**(is ?L ≤ ?R)**  
 <proof>

end

end

## 8 Cutoff Level

This section verifies that the cutoff will be below  $q-max$  with high probability. The result will be needed in Section 9, where it is shown that the estimates will be accurate for any cutoff below  $q-max$ .

**theory** *Distributed-Distinct-Elements-Cutoff-Level*

**imports**

*Distributed-Distinct-Elements-Accuracy-Without-Cutoff*

*Distributed-Distinct-Elements-Tail-Bounds*

**begin**

**hide-const** (**open**) *Quantum.Z*

**unbundle** *intro-cong-syntax*

**lemma** *mono-real-of-int: mono real-of-int*

<proof>

**lemma** *Max-le-Sum:*

**fixes**  $f :: 'a \Rightarrow int$

**assumes** *finite A*

**assumes**  $\bigwedge a. a \in A \implies f a \geq 0$

**shows**  $Max (insert 0 (f ` A)) \leq (\sum a \in A . f a)$  **(is ?L ≤ ?R)**

<proof>

**context** *inner-algorithm-fix-A*

**begin**

The following inequality is true for base e on the entire domain ( $x > 0$ ). It is shown in *ln-add-one-self-le-self*. In the following it is established for base 2, where it holds for  $x \geq 1$ .

**lemma** *log-2-estimate:*

**assumes**  $x \geq (1 :: real)$

**shows**  $log 2 (1+x) \leq x$

<proof>

**lemma** *cutoff-eq-7:*

$real X * 2^{powr (-real q-max)} / b \leq 1$

<proof>

**lemma** *cutoff-eq-6:*

**fixes**  $k$

**assumes**  $a \in A$

**shows**  $(\int f. real-of-int (max 0 (int (f a) - int k)) \partial \Psi_1) \leq 2^{powr (-real k)}$  **(is ?L ≤ ?R)**

<proof>

**lemma** *cutoff-eq-5:*

**assumes**  $x \geq (-1 :: real)$

**shows**  $real-of-int \lfloor log 2 (x+2) \rfloor \leq (real c+2) + max (x - 2^c) 0$  **(is ?L ≤ ?R)**

*<proof>*

**lemma** *cutoff-level*:

*measure*  $\Omega \{ \omega. q \ \omega \ A > q\text{-max} \} \leq \delta/2$  (**is**  $?L \leq ?R$ )

*<proof>*

**end**

**unbundle** *no intro-cong-syntax*

**end**

## 9 Accuracy with cutoff

This section verifies that each of the  $l$  estimate have the required accuracy with high probability assuming as long as the cutoff is below  $q\text{-max}$ , generalizing the result from Section 7.

**theory** *Distributed-Distinct-Elements-Accuracy*

**imports**

*Distributed-Distinct-Elements-Accuracy-Without-Cutoff*

*Distributed-Distinct-Elements-Cutoff-Level*

**begin**

**unbundle** *intro-cong-syntax*

**lemma** (**in** *semilattice-set*) *Union*:

**assumes** *finite*  $I \ I \neq \{\}$

**assumes**  $\bigwedge i. i \in I \implies \text{finite } (Z \ i)$

**assumes**  $\bigwedge i. i \in I \implies Z \ i \neq \{\}$

**shows**  $F \ (\bigcup (Z \ ' I)) = F \ ((\lambda i. (F \ (Z \ i))) \ ' I)$

*<proof>*

This is similar to the existing *hom-Max-commute* with the crucial difference that it works even if the function is a homomorphism between distinct lattices. An example application is  $\text{Max} \ (\text{int} \ ' A) = \text{int} \ (\text{Max} \ A)$ .

**lemma** *hom-Max-commute'*:

**assumes** *finite*  $A \ A \neq \{\}$

**assumes**  $\bigwedge x \ y. x \in A \implies y \in A \implies \text{max} \ (f \ x) \ (f \ y) = f \ (\text{max} \ x \ y)$

**shows**  $\text{Max} \ (f \ ' A) = f \ (\text{Max} \ A)$

*<proof>*

**context** *inner-algorithm-fix-A*

**begin**

**definition**  $t_c$

**where**  $t_c \ \psi \ \sigma = (\text{Max} \ ((\lambda j. \tau_1 \ \psi \ A \ \sigma \ j + \sigma) \ ' \{..<b\})) - b\text{-exp} + 9$

**definition**  $s_c$

**where**  $s_c \ \psi \ \sigma = \text{nat} \ (t_c \ \psi \ \sigma)$

**definition**  $p_c$

**where**  $p_c \ \psi \ \sigma = \text{card} \ \{j \in \{..<b\}. \tau_1 \ \psi \ A \ \sigma \ j + \sigma \geq s_c \ \psi \ \sigma\}$

**definition**  $Y_c$

**where**  $Y_c \ \psi \ \sigma = 2 \wedge s_c \ \psi \ \sigma * \varrho\text{-inv} \ (p_c \ \psi \ \sigma)$

**lemma**  $s_c$ -eq-s:

**assumes**  $(f,g,h) \in \text{sample-pro } \Psi$

**assumes**  $\sigma \leq s f$

**shows**  $s_c(f,g,h) \sigma = s f$

$\langle \text{proof} \rangle$

**lemma**  $p_c$ -eq-p:

**assumes**  $(f,g,h) \in \text{sample-pro } \Psi$

**assumes**  $\sigma \leq s f$

**shows**  $p_c(f,g,h) \sigma = p(f,g,h)$

$\langle \text{proof} \rangle$

**lemma**  $Y_c$ -eq-Y:

**assumes**  $(f,g,h) \in \text{sample-pro } \Psi$

**assumes**  $\sigma \leq s f$

**shows**  $Y_c(f,g,h) \sigma = Y(f,g,h)$

$\langle \text{proof} \rangle$

**lemma** *accuracy-single: measure*  $\Psi \{ \psi. \exists \sigma \leq q\text{-max}. |Y_c \psi \sigma - \text{real } X| > \varepsilon * X \} \leq 1/2^4$

(**is**  $?L \leq ?R$ )

$\langle \text{proof} \rangle$

**lemma** *estimate1-eq:*

**assumes**  $j < l$

**shows** *estimate1*  $(\tau_2 \omega A \sigma, \sigma) j = Y_c(\omega j) \sigma$  (**is**  $?L = ?R$ )

$\langle \text{proof} \rangle$

**lemma** *estimate-result-1:*

*measure*  $\Omega \{ \omega. (\exists \sigma \leq q\text{-max}. \varepsilon * X < |\text{estimate}(\tau_2 \omega A \sigma, \sigma) - X|) \} \leq \delta/2$  (**is**  $?L \leq ?R$ )

$\langle \text{proof} \rangle$

**theorem** *estimate-result:*

*measure*  $\Omega \{ \omega. |\text{estimate}(\tau \omega A) - X| > \varepsilon * X \} \leq \delta$

(**is**  $?L \leq ?R$ )

$\langle \text{proof} \rangle$

**end**

**lemma** (**in** *inner-algorithm*) *estimate-result:*

**assumes**  $A \subseteq \{..<n\} A \neq \{ \}$

**shows** *measure*  $\Omega \{ \omega. |\text{estimate}(\tau \omega A) - \text{real}(\text{card } A)| > \varepsilon * \text{real}(\text{card } A) \} \leq \delta$  (**is**  $?L \leq ?R$ )

$\langle \text{proof} \rangle$

**unbundle** *no intro-cong-syntax*

**end**

## 10 Outer Algorithm

This section introduces the final solution with optimal size space usage. Internally it relies on the inner algorithm described in Section 6, depending on the parameters  $n$ ,  $\varepsilon$  and  $\delta$  it either uses the inner algorithm directly or if  $\varepsilon^{-1}$  is larger than  $\ln n$  it runs  $\frac{\varepsilon^{-1}}{\ln \ln n}$  copies of the inner algorithm (with the modified failure probability  $\frac{1}{\ln n}$ ) using an expander to select its seeds. The theorems below verify that the probability that the relative accuracy of the median of the copies is too large is below  $\varepsilon$ .

**theory** *Distributed-Distinct-Elements-Outer-Algorithm*

**imports**

*Distributed-Distinct-Elements-Accuracy*  
*Prefix-Free-Code-Combinators.Prefix-Free-Code-Combinators*  
*Frequency-Moments.Landau-Ext*  
*Landau-Symbols.Landau-More*

**begin**

**unbundle** *intro-cong-syntax*

The following are non-asymptotic hard bounds on the space usage for the sketches and seeds respectively. The end of this section contains a proof that the sum is asymptotically in  $\mathcal{O}(\ln(\varepsilon^{-1})\delta^{-1} + \ln n)$ .

**definition** *state-space-usage* =  $(\lambda(n, \varepsilon, \delta). 2^{40} * (\ln(1/\delta) + 1) / \varepsilon^2 + \log 2 (\log 2 n + 3))$

**definition** *seed-space-usage* =  $(\lambda(n, \varepsilon, \delta). 2^{30} + 2^{23} * \ln n + 48 * (\log 2(1/\varepsilon) + 16)^2 + 336 * \ln(1/\delta))$

**locale** *outer-algorithm* =

**fixes**  $n :: \text{nat}$

**fixes**  $\delta :: \text{real}$

**fixes**  $\varepsilon :: \text{real}$

**assumes** *n-gt-0*:  $n > 0$

**assumes** *delta-gt-0*:  $\delta > 0$  **and** *delta-lt-1*:  $\delta < 1$

**assumes** *epsilon-gt-0*:  $\varepsilon > 0$  **and** *epsilon-lt-1*:  $\varepsilon < 1$

**begin**

**definition** *n0* **where**  $n_0 = \max(\text{real } n)$  (*exp (exp 5)*)

**definition** *stage-two* **where**  $\text{stage-two} = (\delta < (1/\ln n_0))$

**definition**  $\delta_i :: \text{real}$  **where**  $\delta_i = (\text{if stage-two then } (1/\ln n_0) \text{ else } \delta)$

**definition**  $m :: \text{nat}$  **where**  $m = (\text{if stage-two then } \text{nat } \lceil 4 * \ln(1/\delta) / \ln(\ln n_0) \rceil \text{ else } 1)$

**definition**  $\alpha$  **where**  $\alpha = (\text{if stage-two then } (1/\ln n_0) \text{ else } 1)$

**lemma** *m-lbound*:

**assumes** *stage-two*

**shows**  $m \geq 4 * \ln(1/\delta) / \ln(\ln n_0)$

*<proof>*

**lemma** *n-lbound*:

$n_0 \geq \text{exp (exp 5)}$   $\ln n_0 \geq \text{exp 5}$   $5 \leq \ln(\ln n_0)$   $\ln n_0 > 1$   $n_0 > 1$

*<proof>*

**lemma** *delta1-gt-0*:  $0 < \delta_i$

*<proof>*

**lemma** *delta1-lt-1*:  $\delta_i < 1$

*<proof>*

**lemma** *m-gt-0-aux*:

**assumes** *stage-two*

**shows**  $1 \leq \ln(1/\delta) / \ln(\ln n_0)$

*<proof>*

**lemma** *m-gt-0*:  $m > 0$

*<proof>*

**lemma** *alpha-gt-0*:  $\alpha > 0$

*<proof>*

**lemma** *alpha-le-1*:  $\alpha \leq 1$

*<proof>*



**sublocale**  $I$ : *inner-algorithm*  $n \delta_i \varepsilon$   
 $\langle \text{proof} \rangle$

**abbreviation**  $\Theta$  **where**  $\Theta \equiv \mathcal{E} \ m \ \alpha \ I.\Omega$

**lemma**  $\Theta$ :  $m > 0 \ \alpha > 0 \ \langle \text{proof} \rangle$

**type-synonym** *state* = *inner-algorithm.state list*

**fun** *single* :: *nat*  $\Rightarrow$  *nat*  $\Rightarrow$  *state* **where**  
*single*  $\vartheta \ x = \text{map } (\lambda j. I.\text{single } (\text{pro-select } \Theta \ \vartheta \ j) \ x) \ [0..<m]$

**fun** *merge* :: *state*  $\Rightarrow$  *state*  $\Rightarrow$  *state* **where**  
*merge*  $x \ y = \text{map } (\lambda(x,y). I.\text{merge } x \ y) \ (\text{zip } x \ y)$

**fun** *estimate* :: *state*  $\Rightarrow$  *real* **where**  
*estimate*  $x = \text{median } m \ (\lambda i. I.\text{estimate } (x \ ! \ i))$

**definition**  $\nu$  :: *nat*  $\Rightarrow$  *nat set*  $\Rightarrow$  *state*  
**where**  $\nu \ \vartheta \ A = \text{map } (\lambda i. I.\tau \ (\text{pro-select } \Theta \ \vartheta \ i) \ A) \ [0..<m]$

The following three theorems verify the correctness of the algorithm. The term  $\tau$  is a mathematical description of the sketch for a given subset, while *local.single*, *local.merge* are the actual functions that compute the sketches.

**theorem** *merge-result*:  $\text{merge } (\nu \ \omega \ A) \ (\nu \ \omega \ B) = \nu \ \omega \ (A \cup B)$  (**is**  $?L = ?R$ )  
 $\langle \text{proof} \rangle$

**theorem** *single-result*:  $\text{single } \omega \ x = \nu \ \omega \ \{x\}$  (**is**  $?L = ?R$ )  
 $\langle \text{proof} \rangle$

**theorem** *estimate-result*:  
**assumes**  $A \subseteq \{..<n\} \ A \neq \{\}$   
**defines**  $p \equiv (\text{pmf-of-set } \{..<\text{pro-size } \Theta\})$   
**shows**  $\text{measure } p \ \{\omega. |\text{estimate } (\nu \ \omega \ A) - \text{real } (\text{card } A)| > \varepsilon * \text{real } (\text{card } A)\} \leq \delta$  (**is**  $?L \leq ?R$ )  
 $\langle \text{proof} \rangle$

The function *encode-state* can represent states as bit strings. This enables verification of the space usage.

**definition** *encode-state*  
**where**  $\text{encode-state} = \text{Lf}_e \ I.\text{encode-state } m$

**lemma** *encode-state*: *is-encoding encode-state*  
 $\langle \text{proof} \rangle$

**lemma** *state-bit-count*:  
 $\text{bit-count } (\text{encode-state } (\nu \ \omega \ A)) \leq \text{state-space-usage } (\text{real } n, \varepsilon, \delta)$   
(**is**  $?L \leq ?R$ )  
 $\langle \text{proof} \rangle$

Encoding function for the seeds which are just natural numbers smaller than *pro-size*  $\Theta$ .

**definition** *encode-seed*  
**where**  $\text{encode-seed} = \text{Nb}_e \ (\text{pro-size } \Theta)$

**lemma** *encode-seed*:  
*is-encoding encode-seed*  
 $\langle \text{proof} \rangle$

**lemma** *random-bit-count*:

**assumes**  $\omega < \text{pro-size } \Theta$

**shows**  $\text{bit-count } (\text{encode-seed } \omega) \leq \text{seed-space-usage } (\text{real } n, \varepsilon, \delta)$

(**is** ?L ≤ ?R)

⟨proof⟩

The following is an alternative form expressing the correctness and space usage theorems. If  $x$  is expression formed by *local.single* and *local.merge* operations. Then  $x$  requires *state-space-usage* (*real*  $n, \varepsilon, \delta$ ) bits to encode and *estimate*  $x$  approximates the count of the distinct universe elements in the expression.

For example:

*estimate* (*local.merge* (*local.single*  $\omega$  1) (*local.merge* (*local.single*  $\omega$  5) (*local.single*  $\omega$  1))) approximates the cardinality of  $\{1, 5, 1\}$  i.e. 2.

**datatype** *sketch-tree* = *Single* nat | *Merge* *sketch-tree* *sketch-tree*

**fun** *eval* :: nat ⇒ *sketch-tree* ⇒ *state*

**where**

*eval*  $\omega$  (*Single*  $x$ ) = *single*  $\omega$   $x$  |

*eval*  $\omega$  (*Merge*  $x$   $y$ ) = *merge* (*eval*  $\omega$   $x$ ) (*eval*  $\omega$   $y$ )

**fun** *sketch-tree-set* :: *sketch-tree* ⇒ nat *set*

**where**

*sketch-tree-set* (*Single*  $x$ ) =  $\{x\}$  |

*sketch-tree-set* (*Merge*  $x$   $y$ ) = *sketch-tree-set*  $x$  ∪ *sketch-tree-set*  $y$

**theorem** *correctness*:

**fixes**  $X$

**assumes** *sketch-tree-set*  $t \subseteq \{..<n\}$

**defines**  $p \equiv \text{pmf-of-set } \{..<\text{pro-size } \Theta\}$

**defines**  $X \equiv \text{real } (\text{card } (\text{sketch-tree-set } t))$

**shows**  $\text{measure } p \{ \omega. |\text{estimate } (\text{eval } \omega t) - X| > \varepsilon * X \} \leq \delta$  (**is** ?L ≤ ?R)

⟨proof⟩

**theorem** *space-usage*:

**assumes**  $\omega < \text{pro-size } \Theta$

**shows**

$\text{bit-count } (\text{encode-state } (\text{eval } \omega t)) \leq \text{state-space-usage } (\text{real } n, \varepsilon, \delta)$  (**is** ?A)

$\text{bit-count } (\text{encode-seed } \omega) \leq \text{seed-space-usage } (\text{real } n, \varepsilon, \delta)$  (**is** ?B)

⟨proof⟩

**end**

The functions *state-space-usage* and *seed-space-usage* are exact bounds on the space usage for the state and the seed. The following establishes asymptotic bounds with respect to the limit  $n, \delta^{-1}, \varepsilon^{-1} \rightarrow \infty$ .

**context**

**begin**

Some local notation to ease proofs about the asymptotic space usage of the algorithm:

**private definition** *n-of* :: *real* × *real* × *real* ⇒ *real* **where** *n-of* =  $(\lambda(n, \varepsilon, \delta). n)$

**private definition** *δ-of* :: *real* × *real* × *real* ⇒ *real* **where** *δ-of* =  $(\lambda(n, \varepsilon, \delta). \delta)$

**private definition** *ε-of* :: *real* × *real* × *real* ⇒ *real* **where** *ε-of* =  $(\lambda(n, \varepsilon, \delta). \varepsilon)$

**private abbreviation** *F* :: (*real* × *real* × *real*) *filter*

**where**  $F \equiv (\text{at-top} \times_F \text{at-right } 0 \times_F \text{at-right } 0)$

**private lemma** *var-simps*:

$n\text{-of} = \text{fst}$

$\varepsilon\text{-of} = (\lambda x. \text{fst} (\text{snd } x))$

$\delta\text{-of} = (\lambda x. \text{snd} (\text{snd } x))$

$\langle \text{proof} \rangle$  **lemma** *evt-n*: eventually  $(\lambda x. n\text{-of } x \geq n)$   $F$

$\langle \text{proof} \rangle$  **lemma** *evt-n-1*:  $\forall_F x \text{ in } F. 0 \leq \text{ln } (n\text{-of } x)$

$\langle \text{proof} \rangle$  **lemma** *evt-n-2*:  $\forall_F x \text{ in } F. 0 \leq \text{ln } (\text{ln } (n\text{-of } x))$

$\langle \text{proof} \rangle$  **lemma** *evt- $\varepsilon$* : eventually  $(\lambda x. 1/\varepsilon\text{-of } x \geq \varepsilon \wedge \varepsilon\text{-of } x > 0)$   $F$

$\langle \text{proof} \rangle$  **lemma** *evt- $\delta$* : eventually  $(\lambda x. 1/\delta\text{-of } x \geq \delta \wedge \delta\text{-of } x > 0)$   $F$

$\langle \text{proof} \rangle$  **lemma** *evt- $\delta$ -1*:  $\forall_F x \text{ in } F. 0 \leq \text{ln } (1 / \delta\text{-of } x)$

$\langle \text{proof} \rangle$

**theorem** *asymptotic-state-space-complexity*:

$\text{state-space-usage} \in O[F](\lambda(n, \varepsilon, \delta). \text{ln } (1/\delta)/\varepsilon^{\wedge 2} + \text{ln } (\text{ln } n))$

(**is** -  $\in O[?F](?rhs)$ )

$\langle \text{proof} \rangle$

**theorem** *asymptotic-seed-space-complexity*:

$\text{seed-space-usage} \in O[F](\lambda(n, \varepsilon, \delta). \text{ln } (1/\delta) + \text{ln } (1/\varepsilon)^{\wedge 2} + \text{ln } n)$

(**is** -  $\in O[?F](?rhs)$ )

$\langle \text{proof} \rangle$

**definition** *space-usage*  $x = \text{state-space-usage } x + \text{seed-space-usage } x$

**theorem** *asymptotic-space-complexity*:

$\text{space-usage} \in O[\text{at-top} \times_F \text{at-right } 0 \times_F \text{at-right } 0](\lambda(n, \varepsilon, \delta). \text{ln } (1/\delta)/\varepsilon^{\wedge 2} + \text{ln } n)$

$\langle \text{proof} \rangle$

**end**

**unbundle** *no intro-cong-syntax*

**end**

## References

- [1] N. Alon, Y. Matias, and M. Szegedy. The space complexity of approximating the frequency moments. *Journal of Computer and System Sciences*, 58(1):137–147, 1999.
- [2] Z. Bar-Yossef, T. S. Jayram, R. Kumar, D. Sivakumar, and L. Trevisan. Counting distinct elements in a data stream. In *Randomization and Approximation Techniques in Computer Science*, pages 1–10. Springer Berlin Heidelberg, 2002.
- [3] J. Błasiok. Optimal streaming and tracking distinct elements with high probability. *ACM Trans. Algorithms*, 16(1):3:1–3:28, 2020.
- [4] P. Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. *Journal of Computer and System Sciences*, 31(2):182–209, 1985.
- [5] P. B. Gibbons and S. Tirthapura. Estimating simple functions on the union of data streams. In *Proceedings of the Thirteenth Annual ACM Symposium on Parallel Algorithms and Architectures*, SPAA '01, pages 281–291, 2001.
- [6] V. Guruswami, C. Umans, and S. Vadhan. Unbalanced expanders and randomness extractors from parvaresh–vardy codes. *J. ACM*, 56(4), jul 2009.
- [7] D. M. Kane, J. Nelson, and D. P. Woodruff. An optimal algorithm for the distinct elements problem. In *Proceedings of the Twenty-Ninth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, PODS '10, pages 41–52, New York, 2010.
- [8] E. Karayel. Finite fields. *Archive of Formal Proofs*, June 2022. [https://isa-afp.org/entries/Finite\\_Fields.html](https://isa-afp.org/entries/Finite_Fields.html), Formal proof development.

- [9] E. Karayel. Formalization of randomized approximation algorithms for frequency moments. *Archive of Formal Proofs*, April 2022. [https://isa-afp.org/entries/Frequency\\_Moments.html](https://isa-afp.org/entries/Frequency_Moments.html), Formal proof development.
- [10] E. Karayel. Expander graphs. *Archive of Formal Proofs*, March 2023. [https://isa-afp.org/entries/Expander\\_Graphs.html](https://isa-afp.org/entries/Expander_Graphs.html), Formal proof development.
- [11] D. Woodruff. Optimal space lower bounds for all frequency moments. In *Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '04, pages 167–175, USA, 2004. Society for Industrial and Applied Mathematics.