

UNIQUE ENTITY ESTIMATION WITH APPLICATION TO THE SYRIAN CONFLICT

BY BEIDI CHEN, ANSHUMALI SHRIVASTAVA, AND REBECCA C. STEORTS

Entity resolution identifies and removes duplicate entities in large, noisy databases and has grown in both usage and new developments as a result of increased data availability. Nevertheless, entity resolution has tradeoffs regarding assumptions of the data generation process, error rates, and computational scalability that make it a difficult task for real applications. In this paper, we focus on a related problem of unique entity estimation, which is the task of estimating the unique number of entities and associated standard errors in a data set with duplicate entities. Unique entity estimation shares many fundamental challenges of entity resolution, namely, that the computational cost of all-to-all entity comparisons is intractable for large databases. To circumvent this computational barrier, we propose an efficient (near-linear time) estimation algorithm based on locality sensitive hashing. Our estimator, under realistic assumptions, is unbiased and has provably low variance compared to existing random sampling based approaches. In addition, we empirically show its superiority over the state-of-the-art estimators on three real applications. The motivation for our work is to derive an accurate estimate of the documented, identifiable deaths in the ongoing Syrian conflict. Our methodology, when applied to the Syrian data set, provides an estimate of $191,874 \pm 1772$ documented, identifiable deaths, which is very close to the Human Rights Data Analysis Group (HRDAG) estimate of 191,369. Our work provides an example of challenges and efforts involved in solving a real, noisy challenging problem where modeling assumptions may not hold.

1. Introduction. Our work is motivated by a real estimation problem associated with the ongoing conflict in Syria. While deaths are tremendously well documented, it is hard to know how many unique individuals have been killed from conflict-related violence in Syria. Since March 2011, increasing reports of deaths have appeared in both the national and international news. There are many inconsistencies from various media sources, which is inherent due to the data collection process and the fact that reported victims are documented by multiple sources. Thus, our ultimate goal is to determine an accurate number of documented, identifiable deaths (with associated standard errors) because such information may contribute to future transitional justice and accountability measures. For instance, statistical estimates of death counts have been introduced as evidence in national court cases and

Keywords and phrases: Syrian conflict, unique entity estimation, entity resolution, clustering, dimension reduction

42 international tribunals investigating the responsibility of state leaders for
 43 crimes against humanity (Grillo, 2016).

44 The main challenge with reliable death estimation of the Syrian data set
 45 is the fact that individuals who are documented as dead are often dupli-
 46 cated in the data sets. To address this challenge, one could employ entity
 47 resolution (de-duplication or record linkage), which refers to the task of re-
 48 moving duplicated records in noisy datasets that refer to the same entity
 49 (Tancredi and Liseo, 2011; Sadinle et al., 2014; Bhattacharya and Getoor,
 50 2006; Baxter et al., 2003; Gutman, Afendulis and Zaslavsky, 2013; Winkler,
 51 2004; McCallum and Wellner, 2004; Deming and Glasser, 1959; Fellegi and
 52 Sunter, 1969). Entity resolution is fundamental in many large data process-
 53 ing applications. Informally, let us assume that each entity (records) is a
 54 vector in \mathbb{R}^D . Then given a data set of M records aggregated from many
 55 data sources with possibly numerous duplicated entities perturbed by noise,
 56 the task of entity resolution is to identify and remove the duplicate entities.
 57 For a review of entity resolution see (Winkler, 2006; Christen, 2012; Liseo
 58 and Tancredi, 2013).

59 One important subtask of entity resolution is estimating the number of
 60 unique entities (records) n out of $M > n$ duplicated entities, which we call
 61 *unique entity estimation*. Entity resolution is a more difficult problem be-
 62 cause it requires one to link each entity to its associated duplicate entities.
 63 To obtain high-accuracy entity resolution, the algorithms must at least evalu-
 64 ate a significant amount of pairs for potential duplicates to ensure a link is
 65 not missed. Due to this (and to the best of our knowledge), accurate entity
 66 resolution algorithms scale quadratically or higher ($> O(M^2)$) making them
 67 computationally intractable for large data sets. Reducing the computational
 68 cost in entity resolution is known as blocking, which, via deterministic or
 69 probabilistic algorithms, places similar records into blocks or bins (Christen,
 70 2012; Steorts et al., 2014). The computational efficiency comes at the cost
 71 of missed links and reduced accuracy for entity resolution. Further, it is not
 72 clear if we can use these crude but cheap entity resolution sub-routines for
 73 unbiased estimation of unique entities with strong statistical guarantees.

74 The primary focus of this paper is on developing a *unique entity estima-*
 75 *tion* algorithm that is motivated by the ongoing conflict in Syria and has
 76 the following desiderata:

- 77 1. The estimation cost should be significantly less than quadratic ($O(M^2)$).
 78 In particular, any methodology requiring one to evaluate all pairs for
 79 linkage is not suitable. This is crucial for the Syrian data set and other
 80 large, noisy data sets (Section 1.3).
- 81 2. To ensure accountability regarding estimating the unique number of

documented identifiable victims in the Syrian conflict, it is essential to understand the statistical properties of any proposed estimator. Such a requirement eliminates many heuristics and rule-based entity resolution tasks, where the estimates may be very far from the true value.

3. In most real entity resolution tasks, duplicated data can occur with arbitrarily large changes including missing information, which we observe in the Syrian data set, and standard modeling assumptions may not hold due to the noise inherent in the data. Due to this, we prefer not to make strong modeling assumptions regarding the data generation process.

1.1. *Related Work for Unique Entity Estimation.* The three aforementioned desiderata eliminate all but random sampling-based approaches. In this section, we review them briefly.

To our knowledge, only two random sampling based methodologies satisfy such requirements. Frank (1978) proposed sampling a large enough sub-graph to estimate the total number of connected components based on the properties of the sub-sampled subgraph. Also, Chazelle, Rubinfeld and Trevisan (2005) proposed finding connected components with high probability by sampling random vertices and then visiting their associated components using breadth-first search (BFS). One major issue with random sampling is that most sampled pairs are unlikely to be matches (no edge) providing nearly no information, as the underlying graph is generally very sparse in practice. Randomly sampling vertices and running BFS required by Chazelle, Rubinfeld and Trevisan (2005) are very likely to result in singleton vertices because many records are themselves unique in entity resolution data sets. In addition, finding all possible connections of a given vertex would require $O(M)$ query for edges. A query for edges corresponds to the query for actual link between two records. Sub-sampling a sub-graph, as in Frank (1978), of size s requires $O(s^2)$ edge queries to completely observe it. Thus, s should be reasonably small in order to scale. Unfortunately, requiring a small s hurts the variance of the estimator. We show that the accuracy of both aforementioned methodologies is similar to the non-adaptive variant of our estimator which has provably large variance. In addition, we show both theoretically and empirically that the methodologies based on random sampling lead to poor estimators.

While some methods have recently been proposed for accurate estimation of unique records, they belong to the Bayesian literature and have difficulty scaling due to the curse of dimensionality with Markov chain Monte Carlo

121 Steorts, Hall and Fienberg (2016); Sadinle et al. (2014); Tancredi and Liseo
122 (2011). The evaluation of the likelihood itself is quadratic. Furthermore,
123 they rely on a strong assumption about the specified generative models for
124 the duplicate records. Given such computational challenges with the current
125 state of the methods in the literature, we take a simple approach, especially
126 given the large and constantly growing data sets that we seek to analyze. We
127 focus on practical methodologies that can easily scale to large data sets with
128 minimal assumptions. Specifically, we propose a unique entity estimation
129 algorithm with sub-quadratic cost, which can be reduced to approximating
130 the number of connected components in a graph with sub-quadratic queries
131 for edges (Section 3.1).

132 The rest of the paper proceeds as follows. Section 1.2 provides our moti-
133 vational application from the Syrian conflict and Section 1.3 remarks on
134 the main challenges of the Syrian data set and our proposed methodology.
135 Section 2.1 provides background on variants of locality sensitive hashing
136 (LSH), which is essential to our proposed methodology. Section 3 provides
137 our proposed methodology for unique entity estimation, which is the first
138 formalism of using efficient adaptive LSH on edges to estimate the con-
139 nected components with sub-quadratic computational time. (An example of
140 our approach is given in section 3.2). More specifically, we draw connections
141 between our methodology and random and adaptive sampling in section 3.3,
142 where we show under realistic assumptions that our estimator is theoretic-
143 ally unbiased and has provably low variance. In addition, in section 3.5,
144 we compare random and adaptive sampling for the Syrian data set, illus-
145 trating the strengths of adaptive sampling. In section 3.6, we introduction
146 the variant of LSH used in our paper. Section 3.7 provides our complete
147 algorithm for unique entity estimation. Section 4 provides evaluations of all
148 the related estimation methods on three real data sets from the music and
149 food industries as well as official statistics. Section 5 reports the documented
150 identifiable number of deaths in the Syrian conflict (with a standard error).

151 1.2. *The Syrian Conflict.* Thanks to Human Rights Data Analysis Group
152 (HRDAG), we have access to four databases from the Syrian conflict which
153 cover roughly the same period, namely March 2011 – April 2014, namely, the
154 Violation Documentation Centre (VDC), Syrian Center for Statistics and
155 Research (CSR-SY), Syrian Network for Human Rights (SNHR), and Syria
156 Shuhada website (SS). Each database lists a different number of recorded
157 victims killed in the Syrian conflict, along with available identifying informa-

tion including full Arabic name, date of death, death location, and gender.¹ 158

Since the above information is collected indirectly, such as through friends 159
and religious leaders, or traditional media resources, it naturally comes with 160
many challenges. The data set has biases, spelling errors, and missing values. 161
In addition, it is well known that there are duplicate entities present in 162
the data sets, making estimation more difficult. The ambiguities in Arabic 163
names make the situation significantly worse as there can be a large textual 164
difference between the full and short names in Arabic. (It is not surprising 165
that the Syrian data set has such biases given that the data is collected in 166
the midst of a conflict). 167

Such ambiguities and lack of additional information make entity resolu- 168
tion on this data set considerably challenging (Price et al., 2014). Owing to 169
the significance of the problem, HRDAG has provided labels for a large sub- 170
set of the data set. More specifically, five different human experts from the 171
HRDAG manually reviewed pairs of records in the four data sets, classifying 172
them as matches if referred to the same entity and non-matches otherwise. 173
Our first goal is to accurately estimate the number of unique victims. Ob- 174
taining a match or non-match label of a given record pair may require mo- 175
mentous cost such as manual human supervision or involving sophisticated 176
machine learning. Given that coming up with hand-matched data is a costly 177
process, *our second goal* is to provide a proxy, automated mechanism to 178
create labeled data. (More information regarding the Syrian data set can be 179
found in Appendix A). 180

1.3. *Challenges and Proposed Solutions.* Consider evaluating the Syr- 181
ian data set using all-to-all records comparisons to remove duplicate enti- 182
ties. With approximately 354,000 records from the Syrian data set, we have 183
around 63 billion pairs (6.3×10^{10}). Therefore, it is impractical to classify all 184
these pairs as matches/non-matches reliably. We cannot expect a few experts 185
(five in our case) to manually label 63 billion pairs. A simple computation of 186
all pairwise similarity (63 billion) takes more than 8 days on a heavyweight 187
machine that can run 56 threads in parallel (28 cores in total). In general, 188
this quadratic computational cost is widely considered infeasible for large 189
data sets. Algorithmic labeling of every pair, even if possible for relatively 190
small datasets, is neither reliable nor efficient. Furthermore, it is hard to 191
understand the statistical properties of algorithmic labelling of pairs. Such 192
challenges, therefore, motivate us to focus on the estimation algorithm with 193
constraints mentioned in Section 1. 194

¹These databases include documented identifiable victims and not those who are miss-
ing in the conflict, hence, any estimate reported only refers to the data at hand.

195 **Our Contributions:** We formalize unique entity estimation as approxi-
 196 mating the number of connected components in a graph with sub-quadratic
 197 $\ll O(M^2)$ computational time. We then propose a generic methodology
 198 that provides an estimate in sample (with standard errors). Our proposal
 199 leverages locality sensitive hashing (LSH) in a novel way for the estima-
 200 tion process, with the required computational complexity that is less than
 201 quadratic. Our proposed estimator is unbiased and has provably low variance
 202 compared to random sampling based approaches. To the best of our knowl-
 203 edge this is the first use of LSH for unique entity estimation in an entity
 204 resolution setting. Our unique entity estimation procedure is broadly appli-
 205 cable to many applications, and we illustrate this on three additional real,
 206 fully labelled, entity resolution data sets, which include the food industry,
 207 the music industry, and an application from official statistics. In the absence
 208 of ground truth information, we estimate that the number of documented
 209 identifiable deaths for the Syrian conflict is 191,874, with standard devia-
 210 tion of 1,772, reported casualties, which is very close to the 2014 HRDAG
 211 estimate of 191,369. This clearly demonstrates the power of our efficient esti-
 212 mator in practice, which does not rely on any strong modeling assumptions.
 213 Out of 63 billion possible pairs, our estimator only queries around 450,000
 214 adaptively sampled pairs ($\simeq O(M)$) for labels, yielding a 99.99% reduction.
 215 The labelling was done using support vector machines (SVMs) trained on a
 216 small number of hand-matched, labeled examples provided by five domain
 217 experts. Our work is an example of the efforts required to solve a real noisy
 218 challenging problem where modeling assumptions may not hold.

219 **2. Variants of Locality Sensitive Hashing (LSH).** In this section,
 220 we first provide a review of LSH and min-wise hashing, which is crucial to our
 221 proposed methodology. We then introduce a variant of LSH — Densified One
 222 Permutation Hashing (DOPH), which is essential to our proposed algorithm
 223 for unique entity estimation in terms of scalability. We first provide a brief
 224 literature review of LSH.

225 *2.1. Review of Locality Sensitive Hashing (LSH).* In this section, we first
 226 provide a review of locality sensitive hashing and min-wise hashing, which
 227 is crucial to our proposed methodology.

228 Locality sensitive hashing (LSH) is a well-known *probabilistic method*
 229 of dimension reduction, which is widely used in computer science and in
 230 database engineering as a way of rapidly finding approximate nearest neigh-
 231 bors (Gionis et al., 1999). More recently, locality sensitive hashing has been
 232 utilized has a form of blocking in entity resolution, where one tries to achieve
 233 scalability and avoid all-to-all record comparisons by placing records into

“partitions” or “blocks” either using deterministic or probabilistic methods. 234

Unlike other conventional forms of dimension reduction or blocking for 235
 entity resolution, LSH uses all the features of a record, and can be ad- 236
 justed to ensure that blocks are manageably small, but then do not allow 237
 for further record linkage within blocks. For example, Vatsalan et al. (2014) 238
 introduced novel data structures for sorting and fast approximate nearest- 239
 neighbor look-up within blocks produced by LSH. Their approach gave a 240
 good balance between speed and recall, but their technique is very spe- 241
 cific to nearest neighbor search. In other related work, Steorts et al. (2014) 242
 proposed clustering-based blocking schemes that are variants on LSH. The 243
 first, transitive locality sensitive hashing (TLSH) is based upon the com- 244
 munity discovery literature such that *a soft transitivity* (or relaxed form of 245
 transitivity) can be imposed across blocks. The second, k -means locality sen- 246
 sitive hashing (KLSH) is based upon the information retrieval literature and 247
 clusters similar records into blocks using a vector-space representation and 248
 projections (KLSH had been used before in information retrieval but never 249
 with entity resolution (Paulevé, Jégou and Amsaleg, 2010)). Steorts et al. 250
 (2014) showed that both KLSH and TLSH gave improvements over popular 251
 methods in the literature such as traditional blocking, canopies (McCallum, 252
 Nigam and Ungar, 2000), and k -nearest neighbors clustering. 253

There are many variants of LSH and one popular form is min-wise hashing. 254
 All LSH methods are defined by a type of similarity and a type of dimension 255
 reduction (Broder, 1997a). Recently, Shrivastava and Li (2014a) showed that 256
 min-wise hashing based approaches are superior to random projection based 257
 approaches when the data is very sparse and feature poor. Furthermore, im- 258
 provements in computational speed can be obtained by using the recently 259
 proposed densification scheme known as densified one permutation hashing 260
 (DOPH) (Shrivastava and Li, 2014a,b). Specifically, the authors proposed 261
 an efficient substitute for min-wise hashing, which only requires one per- 262
 mutation (or one hash function) for generating many different hash values 263
 needed for indexing. In short, the algorithm is linear (or constant) in the 264
 tuning parameters, making it very computationally efficient. 265

2.2. *Shingling.* In entity resolution tasks, each record can be represented 266
 as a string of information. For example, each record in the Syrian data set 267
 can be represented as a short *text* description of the person who died in 268
 the conflict. In this paper, we use a k -grams based shingle representation, 269
 which is the most common representation of text data and naturally gives 270
 a set token (or k -grams). That is, each record is treated as a string and 271
 is replaced by a “bag” (or “multi-set”) of length- k contiguous sub-strings 272

that it contains. Since we will use a k -gram based approach to transform the records, our representation of each record will also be a set, which consists of all the k -contiguous characters occurring in record string. As an illustration, for the record BAKER, TED, we separate it into a 2-gram representation. The resulting set is the following:

BA, AK, KE, ER, RT, TE, ED.

In another example, consider Sammy, Smith, whose 2-gram set representation is

SA, AM, MM, MY, YS, MS, SM, MI, IT, TH.

266 We now have two records that have been transformed into a 2-gram repre-
 267 sentation. Thus, for every record (string) we obtain a set $\subset \mathcal{U}$, where the
 268 universe \mathcal{U} is the set of all possible k -contiguous characters.

269 **2.3. Locality Sensitive Hashing.** In this paper, we leverage LSH, which
 270 comes with sound mathematical formalism and guarantees. LSH is widely
 271 used in computer science and database engineering as a way of rapidly find-
 272 ing approximate nearest neighbors (Indyk and Motwani, 1998; Gionis et al.,
 273 1999). Specifically, the variant of LSH that we utilize is scalable to large
 274 databases, and allows for similarity based sampling of entities in less than a
 275 quadratic amount of time.

276 In LSH, a hash function is defined as $y = h(x)$, where y is the *hash code*
 277 and $h(\cdot)$ the *hash function*. A *hash table* is a data structure that is composed
 278 of *buckets* (not to be confused with blocks), each of which is indexed by a
 279 *hash code*. Each reference item x is placed into a bucket $h(x)$.

280 More precisely, LSH is a family of functions that map vectors to a discrete
 281 set, namely, $h : \mathbb{R}^D \rightarrow \{1, 2, \dots, M\}$, where M is in finite range. Given this
 282 family of functions, similar points (entities) are likely to have the same hash
 283 value compared to dissimilar points (entities). The notion of similarity is
 284 specified by comparing two vectors of points (entities), x and y . We will
 285 denote a general notion of similarity by $\text{SIM}(x, y)$. In this paper, we only
 286 require a relaxed version LSH, and we define this below. Formally, a LSH is
 287 defined by the following definition below:

288 **DEFINITION 1.** (*Locality Sensitive Hashing (LSH)*) Let $x_1, x_2, y_1, y_2 \in$
 289 \mathbb{R}^D and suppose h is chosen uniformly from a family \mathcal{H} . Given a similarity
 290 metric, $\text{SIM}(x, y)$, \mathcal{H} is locality sensitive if $\text{SIM}(x_1, x_2) \geq \text{Sim}(y_2, y_3)$ then
 291 $\Pr_{\mathcal{H}}(h(x_1) = h(x_2)) \geq \Pr_{\mathcal{H}}(h(y_1) = h(y_2))$, where $\Pr_{\mathcal{H}}$ is the probability
 292 over the uniform sampling of h .

The above definition is sufficient condition for a family of functions to be LSH. While many popular LSH families satisfy the aforementioned property, we only require this condition for the work described herein. For a complete review of LSH, we refer to [Rajaraman and Ullman \(2012\)](#).

2.4. *Minhashing.* One of the most popular forms of LSH is minhashing ([Broder, 1997b](#)), which has two key properties — a type of similarity and a type of dimension reduction. The type of similarity used is the Jaccard similarity and the type of dimension reduction is known as the minwise hash, which we now define.

Let $\{0, 1\}^D$ denote the set of all binary D dimensional vectors, while \mathbb{R}^D refers to the set of all D dimensional vectors (of records). Note that records can be represented as a binary vector (or set) representation via shingling, BoW, or combining these two methods. More specifically, given two record sets (or equivalently binary vectors) $x, y \in \{0, 1\}^D$, the Jaccard similarity between $x, y \in \{0, 1\}^D$ is

$$\mathcal{J} = \frac{|x \cap y|}{|x \cup y|},$$

where $|\cdot|$ is the cardinality of the set.

More specifically, the minwise hashing family applies a random permutation π , on the given set S , and stores only the minimum value after the permutation mapping, known as the *minhash*. Formally, the minhash is defined as $h_\pi^{\min}(S) = \min(\pi(S))$, where $h(\cdot)$ is a hash function.

Given two sets S_1 and S_2 , it can be shown by an elementary probability argument that

$$(1) \quad Pr_\pi(h_\pi^{\min}(S_1) = h_\pi^{\min}(S_2)) = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|},$$

where the probability is over uniform sampling of π . It follows from Equation 1 that minhashing is a LSH family for the Jaccard similarity.

Remark: In this paper, we utilize a shingling based approach, and thus, our representation of each record is likely to be very sparse. Moreover, [Shrivastava and Li \(2014c\)](#) showed that minhashing based approaches are superior compared to random projection based approaches for very sparse datasets.

2.4.1. *Densified One Permutation Hashing (DOPH).* LSH has been utilized for more than two-decades, where one can use LSH to reduce the computational cost of entity resolution. More specifically, the main idea is to

319 only match records which have the same hash values, known as blocking or
 320 indexing. One major issue with LSH is that the step of creating blocks (hash
 321 buckets) is expensive because it requires several hash computations (Liang
 322 et al., 2014; Steorts et al., 2014). However, it was recently shown that the
 323 several minwise hashes of data can be computed in data reading time using
 324 the technique of Densified One Permutation Hashing (DOPH). Subsequent
 325 works (Shrivastava and Li, 2014a,b) improved the statistical properties of
 326 DOPH. Wang, Shrivastava and Ryu (2017) illustrated that using DOPH one
 327 can get significant improvements over LSH, which leads to the fastest ap-
 328 proximate near-neighbor search algorithm. In this paper, we use the most
 329 recent variant of DOPH, which is significantly faster in practice compared to
 330 minwise hashing (Shrivastava, 2017). Since we use a shingle based represen-
 331 tation for textual data, DOPH is ideal for our proposed algorithm because
 332 the cost for blocking is the same as the data reading cost, which is about
 333 100 times faster than traditional minwise hashing. Throughout the rest of
 334 the paper, when we refer to minwise hashing will refer to the DOPH algo-
 335 rithm for computing minhashes. Further details of LSH and DOPH can be
 336 found in the aforementioned papers. In addition, we specify another reason
 337 for using LSH as the only blocking mechanism which suits our purpose in
 338 section 3.6.4.

339 **3. Unique Entity Estimation.** In this section, we provide notation
 340 used throughout the rest of the paper and provide an illustrative example.
 341 We then propose our estimator, which is unbiased and has provably low
 342 variance. In addition, random sampling is a special case of our procedure
 343 as explained in section 3.5. Finally, we present our unique entity estimation
 344 algorithm in section 3.3.

3.1. *Notation.* The problem of unique entity estimation can be reduced
 to approximating the number of connected components in a corresponding
 graph. Given a data set with size M , we denote the records as

$$R = \{R_i | 1 \leq i \leq M, i \in \mathbb{Z}\}.$$

Next, we define

$$Q(R_i, R_j) = \begin{cases} 1, & \text{if } R_i, R_j \text{ refer to the same entity.} \\ 0, & \text{otherwise.} \end{cases}$$

Let us represent the data set by a graph $G^* = (E, V)$, with vertices E, V .
 Let vertex V_i correspond to record R_i and vertex V_j correspond to record

R_j . Then let edge E_{ij} represent the linkage between records of R_i and R_j (or vertex V_i and V_j). More specifically, we can represent this by the following relationship:

$$V = \{R_i | 1 \leq i \leq M, i \in \mathbb{Z}\}, \text{ and } E = \{(R_i, R_j) | \forall 1 \leq i, j \leq M, Q(R_i, R_j) = 1\}.$$

3.2. *Illustrative Example.* In this section, we provide an illustrative example of how six records are mapped to a graph G^* . Consider record 3 (John) and record 5 (Johnathan) which correspond to the same entity (John Schaech). In G^* , there is an edge E_{35} that connect these records, denoted by V_3 and V_5 . Now consider records 2, 4, and 6, which all refer to the same entity (Nicholas Cage). In G^* , there are edges E_{24} , E_{26} , and E_{46} that connect these records, denoted by V_2 , V_4 , and V_6 . Observe that each connected component in G^* is a unique entity and also a clique. Therefore, our task is reduced to estimating the number of connected components in G^* .

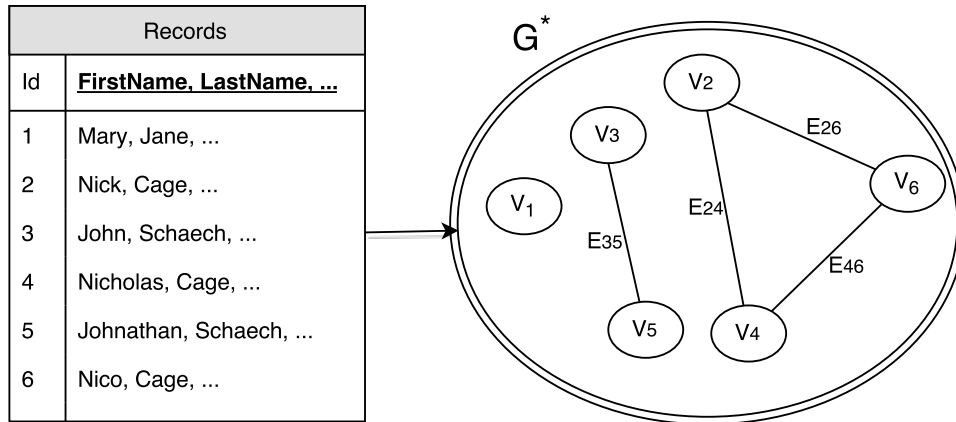


Fig 1: A toy example of mapping records to a graph, where vertices represent records and edges refer to the relation between records.

3.3. *Proposed Unique Entity Estimator.* In this section, we propose our unique entity estimator and provide assumptions that are necessary for our estimation procedure to be practical (scalable).

Since we do not observe the edges of G^* (the linkage), inferring whether there is an edge between two nodes (or whether two records are linked) can be costly, i.e., $O(M^2)$. Hence, one is constrained to probe a small set $\mathcal{S} \subset V \times V$ with $|\mathcal{S}| \ll O(M^2)$ of pairs and query if they have edges. The aim is to use the information about \mathcal{S} to estimate the total number of connected components accurately. More precisely, given the partial graph $G' = \{V, E'\}$,

363 where $E' = E \cap \mathcal{S}$, one wishes to estimate the connected components n of
 364 $G^* = \{V, E\}$.

365 One key property of our estimation process is that we do not make any
 366 modeling assumptions of how duplicate records are generated, and it is not
 367 immediately clear how we can obtain unbiased estimation. For sake of sim-
 368 plicity, we first assume the existence of an efficient (sub-quadratic) process
 369 that samples a small set (near-linear size) of edges \mathcal{S} , such that every edge
 370 in the original graph G^* has (reasonably high) probability p of being in \mathcal{S} .
 371 Thus, set \mathcal{S} , even though small, contains p fraction of the actual edges. For
 372 sparse graphs, as in the case of duplicate records, such a sampler will be
 373 far more efficient than random sampling. Based on this assumption, we will
 374 first describe our estimator and its properties. We then show why our as-
 375 sumption about existence of adaptive sampler is practical by providing an
 376 efficient sampling process based on LSH (Section 3).

377 **Remark:** It is not difficult to see that random sampling is a special case
 378 when $p = \frac{|\mathcal{S}|}{O(M^2)}$ which, as we show later, is a very small number for any
 379 accurate estimation.

380 Our proposed estimator and corresponding algorithm obtains the set of
 381 vertex pairs (or edges) \mathcal{S} through an efficient (adaptive) sampling process
 382 and queries whether there is an edge (linkage) between each pair in \mathcal{S} . Re-
 383 spectively, after the ground truth querying, we observe a sub-sampled graph
 384 G' , consisting of vertices returned by the sampler. Let n'_i be the number of
 385 connected component of size i in the observed graph G' , i.e., n'_1 is the num-
 386 ber of singleton vertices, n'_2 is the number of isolated edges, etc. in G' . It
 387 is worth noting that every connected component in G' is a part of some clique
 388 (maybe larger) in G^* . Let n_i^* denote the number of connected components
 389 (clique) of size i in the original (unobserved) graph G^* .

390 Observe that under the sampling process, any original connected compo-
 391 nent, say C_i^* (clique), will be sub-sampled and can appear as some possibly
 392 smaller connected component in G' . For example, a singleton set in G^* will
 393 remain the same in G' . An isolated edge, on the other hand, can appear as
 394 an edge in G' with probability p and as two singleton vertices in G' with
 395 probability $1 - p$. A triangle can decompose into three possibilities with
 396 probability shown in figure 2. Each of these possibilities provides a linear
 397 equation connecting n_i^* to n'_i . These equations up to cliques of size three are

$$(2) \quad \mathbb{E}[n'_3] = n_3^* \cdot p^2 \cdot (3 - 2p)$$

$$(3) \quad \mathbb{E}[n'_2] = n_2^* \cdot p + n_3^* \cdot (3 \cdot (1 - p)^2 \cdot p)$$

$$(4) \quad \mathbb{E}[n'_1] = n_1^* + n_2^* \cdot (2 \cdot (1 - p)) + n_3^* \cdot (3 \cdot (1 - p)^2).$$

Since we observe n'_i , we can solve for the estimator of each n_i^* and compute the number of connected components by summing up all n_i^* .

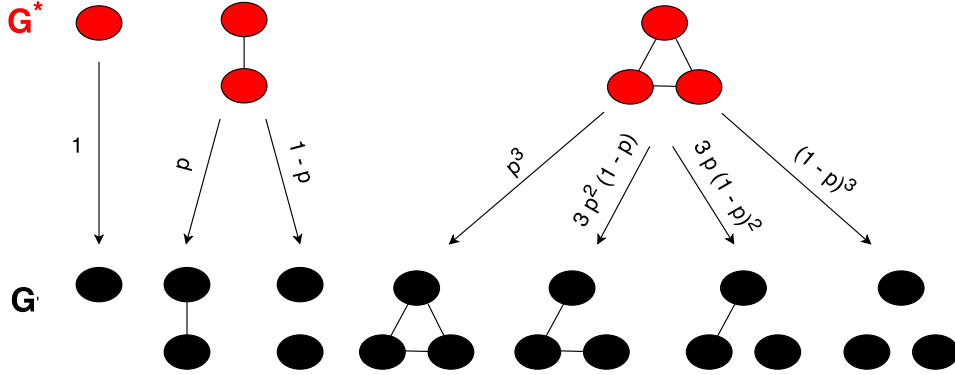


Fig 2: A general example illustrating the transformation and probabilities of connected components from G^* to G' .

Unfortunately, this process quickly becomes combinatorial, and in fact, is at least $\#P$ hard (Provan and Ball, 1983) to compute for cliques of larger sizes. A large clique of size k can appear as many separate connected components and the possibilities of smaller size components it can break into are exponential (Aleksandrov, 1956). Fortunately, we can safely ignore large connected components without significant loss in estimation for two reasons. First, in practical entity resolution tasks, when M is large and contains at least one string-valued feature, it is observed that *most* entities are replicated no more than three or four times. Second, a large clique can only induce large errors if it is broken into many connected components due to undersampling. According to Erdos and Rényi (1960), it will almost surely stay connected if p is high, which is the case with our sampling method.

Assumption: As argued above, we safely assume that the cliques of sizes equal to or larger than 4 in the original graph would retain their structures, i.e., $\forall i \geq 4$, $n_i^* = n'_i$. With this assumption, we can write down the formula for estimating n_1^* , n_2^* , n_3^* by solving Equations 2–4 as,

$$(5) \quad n_3^* = \frac{\mathbb{E}[n'_3]}{p^2 \cdot (3 - 2p)}, \quad n_2^* = \frac{\mathbb{E}[n'_2] - n_3^* \cdot (3 \cdot (1 - p)^2 \cdot p)}{p}$$

$$(6) \quad n_1^* = \mathbb{E}[n'_1] - n_2^* \cdot (2 \cdot (1 - p)) - n_3^* \cdot (3 \cdot (1 - p)^2)$$

It directly follows that our estimator, which we call the Locality Sensitive Hashing Estimator (LSHE) for the number of connected components is given

by

$$(7) \quad \text{LSHE} = n'_1 + n'_2 \cdot \frac{2p-1}{p} + n'_3 \cdot \frac{1-6 \cdot (1-p)^2 \cdot p}{p^2 \cdot (3-2p)} + \sum_{i=4}^M n'_i.$$

412 3.4. *Optimality Properties of LSHE.* We now prove two properties of our
 413 unique entity estimator, namely, that it is unbiased and that it has provably
 414 lower variance than random sampling approaches.

THEOREM 1. *Assuming $\forall i \geq 4$, $n_i^* = n'_i$, we have*

$$(8) \quad \mathbb{E}[\text{LSHE}] = n \quad \text{unbiased}$$

$$(9) \quad \mathbb{V}[\text{LSHE}] = n_3^* \cdot \frac{(p-1)^2 \cdot (3p^2 - p + 1)}{p^2 \cdot (3-2p)} + n_2^* \cdot \frac{(1-p)}{p}$$

415 *The above estimator is unbiased and the variance is given by Equation 9.*

416 Theorem 2 proves the variance of our estimator is monotonically decreas-
 417 ing with p .

418 THEOREM 2. *$\mathbb{V}[\text{LSHE}]$ is monotonically decreasing when p increases in
 419 range $(0, 1]$.*

420 The proof of Theorem 2 directly follows Lemma 1, which is immediately
 421 given.

422 LEMMA 1. *First order derivative of $\mathbb{V}[\text{LSHE}]$ is negative when $p \in (0, 1]$.*

423 Note that when $p = 1$, $\mathbb{V}[\text{LSHE}] = 0$ which means the observed graph G'
 424 is exactly the same as G^* . For detailed proofs of unbiasedness and Lemma
 425 2, see Appendix B.

3.5. *Adaptive Sampling versus Random Sampling.* Before we describe
 our adaptive sampler, we briefly quantify the advantages of an adaptive
 sampling over random sampling for the Syrian data set by computing the
 differences between their variances. Let p be the probability that an edge
 (correct match) is sampled. On the Syrian data set, our proposed sampler,
 described in next section, empirically achieves $p = 0.83$, by reporting around
 450,000 sampled pairs ($O(M)$) out of the 63 billion possibilities ($O(M^2)$).
 Substituting this value of p , the corresponding variance can be calculated
 from Equation 9 as

$$n_3^* \cdot 0.07 + n_2^* \cdot 0.204.$$

Turning to plain random sampling of edges, in order to achieve the same sample size above leads to p as low as $\frac{4.5 \times 10^5}{6.3 \times 10^{10}} \simeq 6.9 \times 10^{-6}$. With such minuscule p , the resulting variance is

$$n_3^* \cdot 6954620166 + n_2^* \cdot 144443.$$

Thus, the variance for random sampling is roughly 7×10^5 times the number of duplicates in the data set and 1×10^{11} the number of triplets in the data set.

In section 4, we illustrate that two other random sampling based algorithms of Chazelle, Rubinfeld and Trevisan (2005) and Frank (1978) also have poor accuracy compared to our proposed estimator. The poor performance of random sampling is not surprising from a theoretical perspective, and illustrates a major weakness empirically for the task of unique entity estimation with sparse graphs, where adaptive sampling is significantly advantageous.

3.6. *The Missing Ingredient: (K,L)-LSH Algorithm.* Our proposed methodology, for unique entity estimation, assumes that we have an efficient algorithm that adaptively samples a set of record pairs, in sub-quadratic time. In this section, we argue that using a variant of LSH (Section 2.1) we can construct such an efficient sampler.

As already noted, we do not make any modeling assumptions on the generation process of the duplicate records. Also, we cannot assume that there is a fixed similarity threshold, because in real datasets duplicates can have arbitrarily large similarity. Instead, we rely on the observation that record pairs with high similarity have a higher chance of being duplicate records. That is, we assume that when two entities R_i and R_j are similar in their attributes, it is more likely that they refer to the same entities (Christen, 2012).² We note that this probabilistic observation is the weakest possible assumption, and almost always true for entity resolution tasks because linking records by a similarity score is one simple way of approaching entity resolution (Christen, 2012; Winkler, 2006; Fellegi and Sunter, 1969).

The similarity between entities (records) naturally gives us a notion of adaptiveness. One simple adaptive approach is to sample records pairs with probability proportional to their similarity. However, as a prerequisite for such sampling, we must compute all the pairwise similarities and associated probability values with every edge. Computing such a pairwise similarity score is a quadratic operation ($O(M^2)$) and is intractable for large datasets.

²The similarity metric that we use to compare sets of record strings is the Jaccard similarity.

458 Fortunately, recent work has shown that (Spring and Shrivastava, 2017a,b;
 459 Luo and Shrivastava, 2017) it is possible to sample pairs adaptively in pro-
 460 portion to the similarity in provably sub-quadratic time using LSH, which
 461 we describe in the next section.

462 3.6.1. *(K,L)-LSH Algorithm and Sub-quadratic Adaptive Sampling.* We
 463 leverage a very recent observation associated with the traditional (K, L)
 464 parameterized LSH algorithm. The (K, L) parameterized LSH algorithm
 465 is a popular similarity search algorithm, which given a query q , retrieves
 466 element x from a preprocessed data set in sub-linear time ($O(KL) \ll M$)
 467 with probability $1 - (1 - \mathcal{J}(q, x)^K)^L$. Here, \mathcal{J} denotes the Jaccard similarity
 468 between the query and the retrieved data vector x . Our proposed method
 469 leverages this (K, L) -parameterized LSH Algorithm, and we briefly describe
 470 the algorithm in this section. For complete details refer to Andoni and Indyk
 471 (2004).

472 Before we proceed, we define hash maps and keys. We use hash maps,
 473 where every integer (or key) is associated with a bucket (or a list) of records.
 474 In a hash map, searching for the bucket corresponding to a key is a constant
 475 time operation. Please refer to algorithms literature (Rajaraman and Ull-
 476 man, 2012) for details on hashing and its computational complexity. Our
 477 algorithm will require several hash maps, L of them, where a record R_i is
 478 associated with a unique bucket in every hash map. The key correspond-
 479 ing to this bucket is determined by minwise hashes of the record R_i . We
 480 encourage readers to refer to Andoni and Indyk (2004) for implementation
 481 details.

482 More precisely, let h_{ij} , $i = \{1, 2, \dots, L\}$ and $j = \{1, 2, \dots, K\}$ be
 483 $K \times L$ minwise hash functions (Equation 1) with each minwise hash func-
 484 tion formed by independently choosing the underlying permutation π . Next,
 485 we construct L meta-hash functions (or the keys) $H_i = \{h_{i,1}, h_{i,2}, \dots, h_{i,K}\}$,
 486 where each of the H_i 's is formed by combining K different minwise hash
 487 functions. For this variant of the algorithm, we need a total of $K \times L$ func-
 488 tions. With such L meta-hash functions, the algorithm has two main phases,
 489 namely the data pre-processing and the sampling pairs phases, which we
 490 outline below.

491 • **Data Preprocessing Phase:** We create L different hash maps (or
 492 hash tables), where every hash values maps to a bucket of elements. For
 493 every record R_j in the dataset, we insert R_j in the bucket associated
 494 with the key $H_i(R_j)$, in hash map $i = \{1, 2, \dots, L\}$. To assign K -tuples
 495 H_i (meta-hash) to a number in a fixed range, we use some universal
 496 random mapping function to the desired address range. See Andoni

and Indyk (2004); Wang, Shrivastava and Ryu (2017) for details. 497

- **Sample Pair Reporting:** For every record R_j in the dataset and 498
from each table i , we obtain all the elements in the bucket associated 499
with key $H_i(R_j)$, where $i = \{1, 2, \dots, L\}$. We then take the union 500
of the L buckets obtained from the L hash tables, and denote this 501
(aggregated) set by A . We finally, report pairs of records (R_i, R_j) , 502
where $R \in A$. 503

THEOREM 3. *The (K,L) -LSH Algorithm reports a pair (R_i, R_j) with 504
probability $1 - (1 - \mathcal{J}(R_i, R_j)^K)^L$, where $\mathcal{J}(R_i, R_j)$ is the Jaccard Similarity 505
between record pairs (R_i, R_j) . 506*

Proof: Since all the minwise hashes are independent due to an indepen- 507
dent sampling of permutations, the probability that both R_i and R_j belong 508
to the same bucket in any hash table i is $\mathcal{J}(R_i, R_j)^K$. Note from equa- 509
tion 1, each meta-hash agreement has probability $\mathcal{J}(R_i, R_j)$. Therefore, 510
the probability that pair (R_i, R_j) is missed by all the L tables is precisely 511
 $(1 - \mathcal{J}(R_i, R_j)^K)^L$, and thus, the required probability of successful retrieval 512
is the complement. 513

The probabilistic expression $1 - (1 - \mathcal{J}(R_i, R_j)^K)^L$ is a monotonic func- 514
tion of the underlying similarity $Sim(q, y)$ associated with the LSH. In par- 515
ticular, higher similarity pairs have more chance of being retrieved. Thus, 516
LSH provides the required sampling that is adaptive in similarity and is 517
sub-quadratic in running time. 518

3.6.2. Computational Complexity. The computational complexity for sam- 519
pling with M records is $O(MKL)$. The procedure requires computing KL 520
minwise hashes for each record. This step is followed by adding every record 521
to L hash tables. Finally, for each record, we aggregate L buckets to form 522
sample pairs. The result of monotonicity and adaptivity of the samples ap- 523
plies to any value of K and L . We choose $O(K \times L) \ll O(M)$ such that 524
we are able to get samples in sub-quadratic time. We further tune K and L 525
using cross-validation to limit the size of our samples. In section 5.3, we evalu- 526
ate the effect of varying K and L in terms of the recall and reduction ratio. 527
(For a review of the recall and reduction ratio, we refer to Christen (2012).) 528
We address the precision at the very end of our experimental procedure to 529
ensure that the recall, reduction ratio, and precision of our proposed unique 530
entity estimation procedure are all as close to 1 as possible while ensuring 531
that the entire algorithm is computationally efficient. For example, on the 532
Syrian data set, we can generate 450,000 samples in less than 127 sec with an 533
adaptive sampling probability (recall) p as high as 0.83. (Note: the prepro- 534

535 cessing is of the order of data loading cost using the (K,L)-LSH Algorithm).
 536 On the other hand, computing all pairwise similarities (63 billion) takes
 537 more than 8 days on the same machine with 28 cores capable of running 56
 538 threads in parallel. We refer to Sadosky et al. (2015) regarding specific com-
 539 parisons of traditional and advanced blocking methods. Specifically, figures
 540 1–3 illustrate variants of blocking, which perform extremely poorly on the
 541 Syrian data set for two reasons. The first is that the recall and the precision
 542 are both extremely low for entity resolution to be practical. The second rea-
 543 son is that under further inspection the blocks sizes are too large to manage
 544 for entity resolution problems at scale. Hence, our focus in this paper is one
 545 the variant that we find works the best under standard entity resolution
 546 evaluation metrics. Next, we describe how this LSH sampler is related to
 547 the adaptive sampler described earlier in Section 3.3.

548 *3.6.3. Underlying Assumptions and Connections with p .* Recall that we
 549 can efficiently sample record pairs R_i, R_j with probability $1 - (1 - \mathcal{J}(R_i, R_j)^K)^L$.
 550 Since we are not making any modeling assumptions, we cannot directly link
 551 this probability to p , the probability of sampling the right duplicated pair
 552 (or linked entities) as required by our estimator LSHE. In the absence of
 553 any knowledge, we can get the estimate of p using a small set of labeled
 554 linked pairs \mathcal{L} . Specifically, we we can estimate the value of p by counting
 555 the fraction of matched pairs (true edges) from \mathcal{L} reported by the sampling
 556 process.

557 Note that in practice there is no similarity threshold θ that guarantees
 558 that two record pairs are duplicate records. That is, it is difficult in practice
 559 to know a fixed θ where $\mathcal{J}(R_i, R_j) \geq \theta$ ensures that R_i and R_j are the
 560 same entities. However, the weakest possible and reasonable assumption is
 561 that high similarity pairs (textual similarity of records) should have higher
 562 chances of being duplicate records than lower similarity pairs.

Formally, this assumption implies that there exists a monotonic function
 f of similarity $\mathcal{J}(R_i, R_j)$ such that the probability of any R_i, R_j being a
 duplicate record is given by $f(\mathcal{J}(R_i, R_j))$. Since our sampling probability
 $1 - (1 - \mathcal{J}(R_i, R_j)^K)^L$ is also a monotonic function of $\mathcal{J}(R_i, R_j)$, we can
 also write

$$f(\mathcal{J}(R_i, R_j)) = g(1 - (1 - \mathcal{J}(R_i, R_j)^K)^L),$$

563 where g is f composed with h^{-1} which is the inverse of $h(x) = 1 - (1 - x^K)^L$.
 564 Unfortunately, we do not know the form of f or g .

Instead of deriving g (or f), which requires additional implicit assump-
 tions on the form of the functions, our process estimates p directly. In partic-
 ular, the estimated value of p is a data dependent mean-field approximation

of g , or rather,

$$p = \mathbb{E}[g(1 - (1 - \mathcal{J}(R_i, R_j)^K)^L)].$$

Crucially, our estimation procedure does not require any modeling assumptions regarding the generation process of the duplicate records, which is significant for noisy data sets, where such assumptions typically break.

3.6.4. *Why LSH?* Although there are several rule-based blocking methodologies, LSH is the only one that is also a random adaptive sampler. In particular, consider a rule-based blocking mechanism, for example on the Syrian data set, which might block on the date of death feature. Such blocking could be a very reasonable strategy for finding candidate pairs. Note that it is still very likely that duplicate records can have different dates of death because the information could be different or misrepresented. In addition, such a blocking method is deterministic, and different independent runs of the blocking algorithm will report the same set of pairs. Even if we find reasonable candidates, we cannot up-sample the linked records to get an unbiased estimate. There will be a systematic bias in the estimates, which does not have any reasonable correction. In fact, random sampling to our knowledge is the only known choice in the existing literature for an unbiased estimation procedure; however, as already mentioned, random uninformative sampling is likely to be very inaccurate.

LSH, on the other hand, can also be used as a blocking mechanism (Steorts et al., 2014). It is, however, more than just a blocking scheme; it is a provably adaptive sampler. Due to randomness in the blocking, different runs of sampler lead to different candidates, unlike deterministic blocking. We can also average over multiple runs to even increase the concentration of our estimates. The adaptive sampling view of LSH has come to light very recently (Spring and Shrivastava, 2017a,b; Luo and Shrivastava, 2017). With adaptive sampling, we get much sharper unbiased estimators than the random sampling approach. To our knowledge, this is the first study of LSH sampling for unique entity estimation.

3.7. *Putting it all Together: Scalable Unique Entity Estimation.* We now describe our scalable unique entity estimation algorithm. As mentioned earlier, assume that we have a data set that contains a text representation of the M records. Suppose that we have a reasonably sized, manually labeled training set \mathcal{T} . We will denote the set of sampled pairs of records given by our sampling process as \mathcal{S} . Note, each element of \mathcal{S} is a pair. Then our scalable entity resolution algorithm consists of three main steps, with the total computational complexity $O(ML + KL + |\mathcal{S}| + |\mathcal{T}|)$. In our case, we

601 will always have $|\mathcal{S}| \ll O(M^2)$ and $KL \ll M$ (in fact, L will be a small
 602 constant), which ensures that the total cost is strictly sub-quadratic. The
 603 complete procedure is summarized in Algorithm 1.

- 604 1. **Adaptively Sample Record Pairs ($O(ML)$):** We regard each record
 605 R_i as a short string and replace it by an “n-grams” based representa-
 606 tion. Then one computes $K \times L$ minwise hashes of each corresponding
 607 string. This can be done in a computationally efficient manner us-
 608 ing the DOPH algorithm (Shrivastava, 2017), which is done in data
 609 reading time. Next, once these hashes are obtained, one applies the
 610 sampling algorithm described in section 3 in order to generate a large
 611 enough sample set, which we denote by \mathcal{S} . For each record, the sam-
 612 pling step requires exactly L hash table queries, which are themselves
 613 $O(1)$ memory lookups. Therefore, the computational complexity of this
 614 step is $O(ML + KL)$.
- 615 2. **Query each Sample Pairs:** Given the set of sampled pairs of records
 616 \mathcal{S} from Step 1, for every pair of records in \mathcal{S} , we query whether these
 617 record pairs are a match or non-match. This step requires, $O(|\mathcal{S}|)$,
 618 queries for the true labels. Here, one can use manually labeled data
 619 if it exists. In the absence of manually labeled data, we can also use
 620 a supervised algorithm, such as support vector machines or random
 621 forests, that is trained on the manually labeled set \mathcal{T} (Section 5).

- (a) **Estimate p :** Given the sampled set of record pairs \mathcal{S} , we need to
 know the value of p , the probability that any given correct pair
 is sampled. To do so, we use the fraction of true pairs sampled
 from the labeled training set \mathcal{T} . The sampling probability p can
 be estimated by computing the fraction of the matched pairs of
 training set records \mathcal{T}_{match} appearing in \mathcal{S} . That is, we estimate
 p (unbiasedly) by

$$p = \frac{|\mathcal{T}_{match} \cap \mathcal{S}|}{|\mathcal{T}_{match}|}.$$

622 If T is stored in a dictionary, then this step can be done on the
 623 fly while generating samples. It only costs $O(\mathcal{T})$ extra work to
 624 create the dictionary.

- (b) **Count Different Connected Components in G' ($O(M + |\mathcal{S}|)$):** The resulting matched sampled pairs, after querying every
 625 sample for actual (or inferred) labels, form the edges of G' . We
 626 now have complete information about our sampled graph G' . We
 627 can now traverse G' and count all sizes of connected components
 628 in G' to obtain n'_1, n'_2, n'_3 and so on. Traversing the graph has
 629
 630

Algorithm 1 LSH-Based Unique Entity Estimation Algorithm

-
- 1: **Input:** Records R , Labeled Set \mathcal{T} , Sample Size m
 - 2: **Output:** $LSHE$
 - 3: $\mathcal{S} = LSHSampling(R)$ (Section 3.6.1)
 - 4: Get \mathcal{T}_{match} be the linked pairs (duplicate entities) in \mathcal{T}
 - 5: $p = \frac{|\mathcal{T}_{match} \cap \mathcal{S}|}{|\mathcal{T}_{match}|}$
 - 6: Query every pair in \mathcal{S} for match/mismatch (get actual labels). (Graph G')
 - 7: $n'_1, n'_2, n'_3 \dots n'_M = Traverse(G')$
 - 8: $LSHE = Equation\ 7(p, n'_1, n'_2, n'_3 \dots n'_M)$
-

Fig 3: Overview of our proposed unique entity estimation algorithm.

computational complexity $O(M + |\mathcal{S}|)$ time using Breadth First 631
Search (BFS). 632

3. **Estimate the Number of Connected Components in G^* ($O(1)$):** 633
Given the values of p , n'_1 , n'_2 , and n'_3 we use equation 7 to compute 634
the unique entity estimator LSHE. 635

4. Experiments. We evaluate the effectiveness of our proposed method- 636
ology on the Syrian data set and three additional real data sets, where the 637
Syrian data set is only partially labeled, while the other three data sets 638
are fully labeled. We first perform evaluations and comparisons on the three 639
fully labeled data sets, and then give an estimate of the documented number 640
of identifiable victims for the Syrian data set. 641

- **Restaurant:** The **Restaurant** data set contains 864 restaurant records 642
collected from Fodor’s and Zagat’s restaurant guides.³ There are a 643
total of 112 duplicate records. Attribute information contains name, 644
address, city, and cuisine. 645
- **CD:** The **CD** data set that includes 9,763 CDs randomly extracted 646
from freeDB.⁴ There are a total of 299 duplicate records. Attribute 647
information consists of 106 total features such as artist name, title, 648
genre, among others. 649
- **Voter:** The **Voter** data has been scraped and collected by [Christen](#) 650
(2014) beginning in October 2011. We work with a subset of this data 651
set containing 324,074 records. There are a total of 68,627 duplicate 652
records. Attribute information contains personal information on voters 653

³Originally provided by Sheila Tejada, downloaded from
<http://www.cs.utexas.edu/users/ml/riddle/data.html>.

⁴<https://hpi.de/naumann/projects/repeatability/datasets/cd-datasets.html>.

DBname	Domain	Size	# Matching Pairs	# Attributes	# Entities
Restaurants	Restaurant Guide	864	112	4	752
CD	Music CDs	9,763	299	106	9,508
Voter	Registration Info	324,074	70,359	6	255,447
Syria	Death Records	354,996	N/A	6	N/A

Table 1: We present five important features of the four data sets. **Domain** reflects the variety of the data type we used in the experiments. **Size** is the number of total records respectively. **# Matching Pairs** shows how many pair of records point to the same entity in each data set. **# Attributes** represents the dimensionality of individual record. **# Entities** is the number of unique records.

654 from North Carolina including full name, age, gender, race, ethnicity,
655 address, zip code, birth place, and phone number.

- 656 • **Syria**: The **Syria** data set comprises data from the Syrian conflict,
657 which covers the same time period, namely, March 2011 – April 2014.
658 This data set is not publicly available and was provided by HRDAG.
659 The respective data sets come from the Violation Documentation Cen-
660 tre (VDC), Syrian Center for Statistics and Research (CSR-SY), Syr-
661 ian Network for Human Rights (SNHR), and Syria Shuhada website
662 (SS). Each database lists a different number of recorded victims killed
663 in the Syrian conflict, along with available identifying information in-
664 cluding full Arabic name, date of death, death location, and gender.⁵

665 The above datasets cover a wide spectrum of different varieties observed
666 in practice. For each data set, we report summary information in Table 1.

Id	First Name	Last Name	Gender	Date of Death	Governorate	Location
1	مبنديز	بيحيدي	F	2011-10-23	Homs	قره‌هاقلا عراشقره‌هاقلا عراش
2	مبنديز	بيحيدي	F	2011-10-23	Homs	عراشعراش
3	مبنديز		F	2011-10-25	Homs	ةميدقلا صمحو

Fig 4: We show several death records in Syrian dataset from VDC, which allows for public access to some of the data. All of the three records belong to the same entity, labeled by human experts. Record 1 and 2 are similar in all attributes while Record 1 and 3 are very different. Due to the variation in the data, records that are very similar are likely to be linked as the same entity, however, it is more difficult to make decisions when records show differences, such as record 1 and 3. This illustrates some of the limitations from deterministic blocking methods discussed in Section 3.6.4.

⁵These databases include documented identifiable victims and not those who are missing in the conflict. Hence, any estimate reported only refers to the data at hand.

4.1. *Evaluation Settings.* In this section, we outline our evaluation settings. We denote Algorithm 1 as the LSH Estimator (LSHE). We make comparisons to the non-adaptive variant of our estimator (PRSE), where we use plain random sampling (instead of adaptive sampling). This baseline uses the same procedure as our proposed LSHE, except that the sampling is done uniformly. A comparison with PRSE quantifies the advantages of the proposed adaptive sampling over random sampling. In addition, we implemented the two other known sampling methods, for connected component estimation, proposed in Frank (1978) and Chazelle, Rubinfeld and Trevisan (2005). For convenience, we denote them as Random Sub-Graph based Estimator (RSGE), and BFS on Random Vertex based Estimator (BFSE) respectively. Since the algorithms are based on sampling (adaptive or random), to ensure fairness, we fix a budget m as the number of pairs of vertices considered by the algorithm. Note that any query for an edge is a part of the budget. If the fixed budget is exhausted, then we stop the sampling process and use the corresponding estimate, using all the information available.

We briefly describe the implementation details of the four considered estimators below:

1. **LSHE:** In our proposed algorithm, we use the (K, L) parameterized LSH algorithm to generate samples of record pairs using Algorithm 3, where recall K and L control the resulting sample size (section 5.3). Given K, L as an input to Algorithm 1, we use the sample size as the value of the fixed budget m . Table 2 gives different sample budget sizes (with the corresponding K and L) and corresponding values of p for selected samples in three real data sets.
2. **PRSE:** For a fair comparison, in this algorithm, we randomly sample the same number of record pairs used by LSHE. We then perform the same estimation process as LSHE but instead use $p = \frac{2m}{M(M-1)}$, which corresponds to the random sampling probability to get the same number of samples, which is m .
3. **RSGE (Frank, 1978):** This algorithm requires performing breadth first search (BFS) on each randomly selected vertices. BFS requires knowing all edges (neighbors) of a node for the next step, which requires $M - 1$ edge queries. To ensure the fixed budget m , we end the traversal when the number of distinct edge queries reaches the fixed budget m .
4. **BFSE (Chazelle, Rubinfeld and Trevisan, 2005):** This algorithm samples a subgraph and observes it completely. This requires labeling all the pairs of records in the sampled sub-graph. To ensure same

706 budget m , the sampled sub-graph has approximately $\sqrt{2m}$ vertices.

707 **Remark:** To the best of our knowledge there have been no experimental
 708 evaluations of the two algorithms of Frank (1978) and Chazelle, Rubinfeld
 709 and Trevisan (2005) in the literature. Hence, our results could be of inde-
 710 pendent interest in themselves.

We compute the relative error (RE), calculated as

$$\text{RE} = \frac{|\text{LSHE} - n|}{n},$$

711 for each of the estimators, for different values of the budget m . We plot the
 712 RE for each of the estimators, over a range of values of m , summarizing the
 713 results in figure 5.

714 All the estimators require querying pairs of records compared to labeled
 715 ground truth data for whether they are a match or a non-match. As already
 716 mentioned, in the absence of full labeled ground truth data, we can use a
 717 supervised classifiers such as SVMs as a proxy, assuming at least some small
 718 amount of labeled data exists. By training an SVM, we can use this as a
 719 proxy for labeled data as well. We use such a proxy in the Syrian data set
 720 because we are not able to query every pair of records to determine whether
 721 they are true duplicates or not.

722 We start with the three data sets where fully labelled ground truth data
 723 exists. For LSHE, we compute the estimation accuracy using both the su-
 724 pervised SVM (Section 5) as well as using the fully labelled ground truth
 725 data. The difference in these two numbers quantifies the loss in estimation
 726 accuracy due to the use of the proxy SVM prediction instead of using ground
 727 truth labeled data. In our use of SVMs, we take less than 0.01% of the total
 728 number of the possible record pairs as the training set.

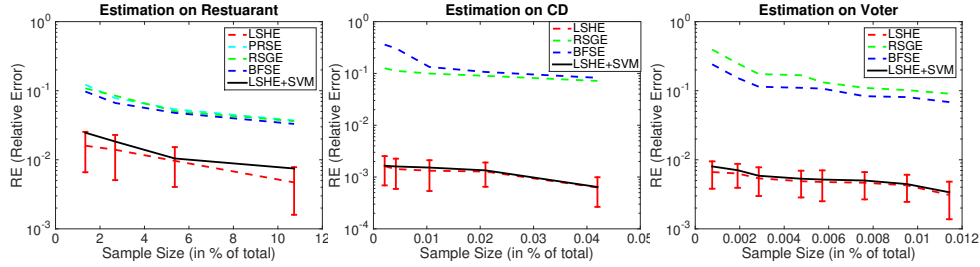


Fig 5: The dashed lines show the RE of the four estimators on the three real data sets, where the y-axis is on the log-scale. Observe that LSHE outperforms all other three estimators in one to two orders of magnitude. The standard deviation of the RE for LSHE is also shown in the plots with the red error bars, which is with respect to randomization of hash functions. In particular, the PRSE performs unreliable estimation on the CD and Voter data sets. The dashed and solid black lines represent RE of LSHE using ground truth labels and a SVM classifier (y-axis is on the log scale). We discuss the LSHE + SVM estimator in section 5 (solid black line).

4.2. *Evaluation Results.* In this section, we summarize our results regarding the aforementioned evaluation metrics by varying the sample size m on the three real data sets (see figure 5). We notice that for the CD and Voter data sets, we cannot obtain any reliable estimate (for any sample size) using PRSE. Recall that plain random sampling almost always samples pairs of records that correspond to non-matches. Thus, it is not surprising that this method is unreliable because sampling random pairs is unlikely to result in a duplicate pair for entity resolution tasks. Even with repeated trials, there are no edges in the specified sampled pairs of records, leading to an undefined value of p . This phenomenon is a common problem in random sampling estimators over sparse graphs. Almost all the sampled nodes are singletons. Subsampling a small sub-graph leads to a graph with most singleton nodes, which leads to a poor accuracy of BFSE. Thus, it is expected that random sampling will perform poorly. Unfortunately, there is no other baseline for unbiased estimation of the number of unique entities.

From figure 5 observe that the RE for proposed estimator LSHE is approximately one to two orders of magnitude lower than the other considered methods, where the y-axis is on the log-scale. Undoubtedly, our proposed estimator LSHE consistently leads to significantly lower RE (lower error rates) than the other three estimators. This is not surprising from the analysis shown in section 3.5. The variance of random sampling based methodologies will be significantly higher.

751 Taking a closer look at LSHE, we notice that we are able to efficiently
 752 generate samples with very high values of p (see Table 2). In addition, we
 753 can clearly see that LSHE achieves high accuracy with very few samples.
 754 For example, for the CD data set, with a sample size less than 0.05% of the
 755 total possible pairs of records of the entire data set, LSHE achieves 0.0006
 756 RE. Similarly, for the Voter data set, with a sample size less than 0.012%
 757 of the total possible pairs of records of the entire data set, LSHE achieves
 758 0.003 RE.

759 Also, note the small values of K and L parameters required to achieve
 760 the corresponding sample size. K and L affect the running time, and small
 761 values $KL \ll O(M^2)$ indicate significant computational savings as argued
 762 in section 3.6.2

763 As mentioned earlier, we also evaluate the effect of using SVM prediction
 764 as a proxy for actual labels with our LSHE. The dotted plot shows those
 765 results. We remark on the results for LSHE + SVM in section 5.

	Restaurant				CD				Voter			
Size	1.0	2.5	5.0	10	0.005	0.01	0.02	0.04	0.002	0.006	0.009	0.013
p	0.42	0.54	0.65	0.82	0.72	0.74	0.82	0.92	0.62	0.72	0.76	0.82
K	1	1	1	1	1	1	1	1	4	4	4	4
L	4	8	12	20	5	6	8	14	25	32	35	40

Table 2: We illustrate part of the sample sizes (in % in TOTAL) for different sets of samples generated by Min-Wise Hashing and their corresponding p in all three data sets.

766 **5. Documented Identifiable Deaths in the Syrian Conflict.** In
 767 this section, we describe how we estimate the number of documented identi-
 768 fiable deaths for the Syrian data set. As noted before, we do not have ground
 769 truth labels for all record pairs, but the data set was partially labeled with
 770 40,000 record pairs (out of 63 billion). We propose an alternative (auto-
 771 matic) method of labeling the sample pairs, which is also needed by our
 772 proposed estimation algorithm. More specifically, using the partially labeled
 773 pairs, we train an SVM. In fact, other supervised methods could be con-
 774 sidered here, such as random forests, Bayesian Adaptive Regression Trees
 775 (BART), among others, however, given that SVMs perform very well, we
 776 omit such comparisons as we expect the results to be similar if not worse.

777 To train the SVM, we take every record pair and generate k -grams rep-
 778 resentation for each record. Then we spilt the partially labeled data into
 779 training and testing sets, respectively. Each training and testing set con-

tains a pair of records $x_k = [R_i, R_j]$. In addition, we can use a binary label 780
 indicating whether the record pair is a match or non-match. That is, we can 781
 write the data as $\{x_k = [R_i, R_j], y_k\}$ as the set difference of the k -grams of 782
 the strings of pairs of records R_i and R_j , respectively. Observe that $y_k = 1$ 783
 if the R_i and R_j is labelled as match and $y_k = -1$ otherwise. Next, we tune 784
 the SVM hyper-parameters using 5-fold cross-validation, and we find the 785
 accuracy of SVM on the testing set was 99.9%. With a precision as high a 786
 0.99, we can reliably query an SVM and now treat this as an expert label. 787

To understand the effect of using SVM prediction as a proxy to label 788
 queries in our proposed unique entity estimation algorithm, we return to 789
 observing the behavior in figure 5. We treat the LSHE estimator on the 790
 other three real datasets as our baseline and compare to LHSE with the 791
 SVM component, where the SVM prediction replaces the querying process 792
 (LSHE +SVM). Observe in figure 5, that the plot for LSH (solid black line) 793
 and LSH+SVM (dotted black line) overlap indicating a negligible loss in per- 794
 formance. This overlap is expected given the high accuracy (high precision) 795
 of the SVM classifier. 796

5.1. *Running Time.* We briefly highlight the speed of the sampling process 797
 since it could be used for on the fly or online unique entity estimation. 798
 The total running time for producing 450,000 sampled pairs (out of a possible 799
 63 billion) used for the LSH sampler (Section 3.6.1) with $K = 15$ and 800
 $L = 10$ is 127 seconds. The preprocessing cost is included in the 127 sec- 801
 onds. The preprocessing is of the order of data loading cost using DOPH. 802
 (For further details on the benchmarking performance of DOPH compared 803
 with other LSH methods, please see Wang, Shrivastava and Ryu (2017)). 804
 On the other hand, it will take approximately take 8 days to compute all 805
 pairwise similarities across the 354,996 Syrian records. Computing the pair- 806
 wise similarities is just the first step for any known adaptive sampling over 807
 pairs based on similarity assuming that we do not use the LSH sampler. 808
 (Note: there are other ways of blocking (Christen, 2012; Sadosky et al., 809
 2015), however as mentioned in Section 3.6.4 they are mostly deterministic 810
 (or rule-based) and do not provide an estimate of the unique entities. 811

5.2. *Unique Number of Documented Identifiable Victims.* In the Syrian 812
 dataset, with 354,996 records and possibly 63 billion (6.3×10^{10}) pairs, our 813
 motivating goal was to estimate the unique number of documented identifi- 814
 able victims. Specifically, in our final estimate, we use 452,728 sampled pairs 815
 that are given by LSHE+SVM ($K = 15, L = 10$) which has approximately 816
 $p = 0.83$ based on the subset of labeled pairs. The sample size was chosen 817
 to balance the computational runtime and the value of p . Specifically, one 818

819 wants high values of p (for a resulting low variance of our estimate) and, to
 820 balance running time, we limit the sample size to be around the total num-
 821 ber of records $O(M)$, to ensure a near linear time algorithm. (Such settings
 822 are determined by the application, but as we have demonstrated they work
 823 for a variety of real entity resolution data sets). We chose the SVM as our
 824 classifier to label the matches and non-matches. The final unique number
 825 of documented identifiable victims in the Syrian data set was estimated to
 826 be $191,874 \pm 1772$, very close to the 191,369 documented identifiable deaths
 827 reported by HRDAG 2014, where their process is described in Appendix A.

828 **5.3. Effects of L , K , on sample size and p .** In this section, we discuss
 829 the sensitivity of our proposed method as we vary the choice of L , K , the
 830 sample size M , and p .

831 We want both $KL \ll M$ as well as the number of samples to be $\ll M^2$,
 832 for the process to be truly sub-quadratic. For accuracy, we want high values
 833 of p , because the variance is monotonic in p , which is also the recall of true
 834 labeled pairs. Thus, there is a natural trade-off. If we sample more, we get
 835 high p but more computations.

836 K and L are the basic parameters of our sampler (Section 3.6.1), which
 837 provide a tradeoff between the computationally complexity and accuracy. A
 838 large value of K makes the buckets sparse exponentially), and thus, fewer
 839 pairs of records are sampled from each table. A large value of L increases
 840 the repetition of hash tables (linearly), which increases the sample size. As
 841 already argued, the computational cost is $O(MKL)$.

842 To understand the behavior of K , L , p , and the computational cost, we
 843 perform a set of experiments on the Syrian dataset. We use n-gram of 2–5,
 844 we vary L from 5–100 by steps of 5 and K takes values 15,18,20,23,25,28,30,32,35.
 845 For all these combinations, we then plot the recall (also the value of p) and
 846 the reduction ratio (RR), which is the percentage of computational savings.
 847 A 99% reduction ratio means that the original space has been reduced to
 848 only having to look at a only 1% of total sampled pairs. Figure 6 shows the
 849 tradeoffs between reduction ratio and recall (or value of p). Every dot in the
 850 figure is one whole experiment.

851 Regardless of the n-gram variation from 2–5, the recall and reduction
 852 ratio (RR) are close to 1 as illustrated in figure 6. We see that an n-gram of
 853 3 overall is most stable in having a recall and RR close to 0.99. We observe
 854 that $K = 15$ and $L = 10$ gives a high recall of around 83% with less than
 855 half a million pairs (out of 63 billion possible) to evaluate ($RR \geq 0.99999$).

856 **6. Discussion.** Motivated by three real entity resolution tasks and the
 857 ongoing Syrian conflict, we have proposed a general, scalable algorithm for

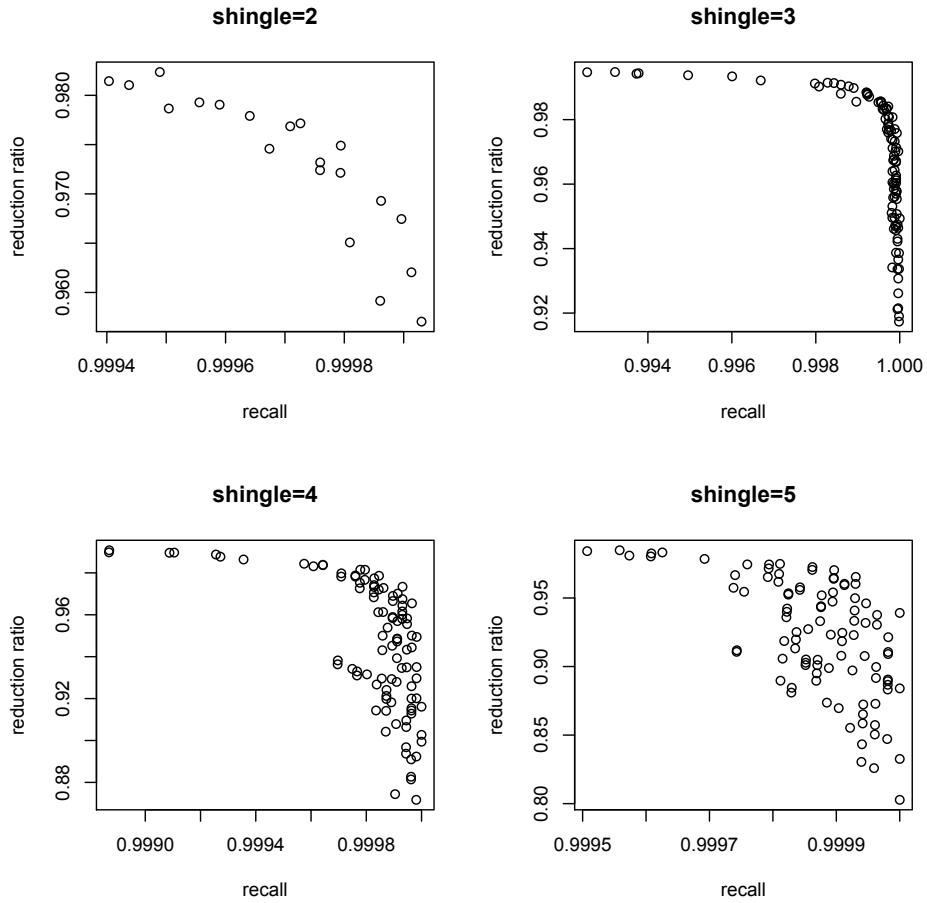


Fig 6: For shingles 2–5, we plot the RR versus the recall. Overall, we see the best behavior for a shingle of 3, where the RR and recall can be reached at 0.98 and 1, respectively. We allow L and K to vary on a grid here. L varies from 5–100 by steps of 5; and K takes values 15, 18, 20, 23, 25, 28, 30, 32, and 35.

858 unique entity estimation. Our proposed method is an adaptive LSH on the
859 edges of a graph, which in turn estimates the connected components in
860 sub-quadratic time. Our estimator is unbiased and has provably low vari-
861 ance in contrast to other such estimators for unique entity estimation in the
862 literature. In experimental results, it outperforms other estimators in the
863 literature on three real entity resolution data sets. Moreover, we have esti-
864 mated the number of documented identifiable deaths to be $191,874 \pm 1772$,
865 which very closely matches the 2014 HRDAG estimate, completed by hand-
866 matching. To our knowledge, we have the first estimate for the number of
867 documented identifiable deaths with a standard error associated with such
868 an estimate. Our methods are scalable, potentially bringing impact to the
869 human rights community, where such estimates could be updated in near
870 real time. It could lead to further impact in public policy and transitional
871 justice in Syria and other areas of conflict globally.

872 **Acknowledgements:** We would like to thank the Human Rights Data
873 Analysis Group (HRDAG) and specifically, Megan Price, Patrick Ball, and
874 Carmel Lee for commenting on our work and giving helpful suggestions that
875 have improved the methodology and writing. We would also like to thank
876 Stephen E. Fienberg and Lars Vilhuber for making this collaboration possi-
877 ble. PhD student Chen is supported by National Science Foundation (NSF)
878 grant number 1652131. Shrivastava's work is supported by NSF-1652131 and
879 NSF-1718478. Steorts's work is supported by NSF-1652431, NSF-1534412,
880 and the Laboratory for Analytic Sciences (LAS). This work is representative
881 of the author's alone and not of the funding organizations.

APPENDIX A: SYRIAN DATA SET

In this section, we provide a more detailed description about the Syrian data set. As mentioned in section 1.2, via collaboration with the Human Rights Data Analysis Group (HRDAG), we have access to four databases. They come from the Violation Documentation Centre (VDC), Syrian Center for Statistics and Research (CSR-SY), Syrian Network for Human Rights (SNHR), and Syria Shuhada website (SS). Each database lists each victim killed in the Syrian conflict, along with identifying information about each person (see Price et al. (2013) for further details).

Data collection by these organizations is carried out in a variety of ways. Three of the groups (VDC, CSR-SY, and SNHR) have trusted networks on the ground in Syria. These networks collect as much information as possible about the victims. For example, information is collected through direct community contacts. Sometimes information comes from a victim’s friends or family members. Other times, information comes from religious leaders, hospitals, or morgue records. These networks also verify information collected via social and traditional media sources. The fourth source, SS, aggregates records from multiple other sources, including NGOs and social and traditional media sources (see <http://syriansshuhada.com/> for information about specific sources).

These lists, despite being products of extremely careful, systematic data collection, are not probabilistic samples (Price, Gohdes and Ball, 2015; Price and Ball, 2015a,b; Price et al., 2014). Thus, these lists cannot be assumed to represent the underlying population of all victims of conflict violence. Records collected by each source are subject to biases, stemming from a number of potential causes, including a group’s relationship within a community, resource availability, and the current security situation. Although it is beyond the scope of this paper, final analyses of these sources must appropriately adjust for such biases before drawing conclusions about patterns of violence.

A.1. Syrian Handmatched Data Set. We describe how HRDAG’s training data on the Syrian data set was created, which we use in our paper. We would like to note that we only use a small fraction of the training data for two reasons. The first is so that we can see how close our estimate is in practice to their original handmatched estimate, given that such a large portion of the data was handmatched. Second, we want to avoid using too much training data to avoid biases and also because such handmatching efforts would not be possible moving forward as the Syrian conflict continues, and our small training data set is meant for one moving forward in practice.

920 First, all documented deaths recorded by any of the documentation groups
921 were concatenated together into a single list. From this list, records were
922 broadly grouped according to governorate and year. In other words, all
923 killings recorded in Homs in 2011 were examined as a group, looking for
924 records with similar names and dates.

925 Next, several experts review these “blocks”, sometimes organized as pairs
926 for comparison and other times organized as entire spreadsheets for review.
927 These experts determine whether pairs or groups of records refer to the same
928 individual victim or not. Pairs or groups of records determined to refer to
929 the same individual are assigned to the same “match group.” All of the
930 records contributing to a single “match group” are then combined into a
931 single record. This new single record is then again examined as a pair or
932 group with other records, in an iterative process.

933 For example, two records with the same name, date, and location may
934 be identified as referring to the same individual, and combined into a single
935 record. In a second review process, it may be found that that record also
936 matches the name and location, but not date, of a third record. The third
937 record may list a date one week later than the two initial records, but still be
938 determined to refer to the same individual. In this second pass, information
939 from this third record will also be included in the single combined record.

940 When records are combined, the most precise information available from
941 each of the individual records is kept. If some records contain contradictory
942 information (for example, if records A and B record the victim as age 19
943 and record C records age 20) the most frequently reported information is
944 used (in this case, age 19). If the same number of records report each piece
945 of contradictory information, a value from the contradictory set is randomly
946 selected.

947 Three of the experts are native Arabic speakers; they review records with
948 the original Arabic content. Two of the experts review records translated
949 into English. These five experts review overlapping sets of records, meaning
950 that some records are evaluated by two, three, four, or all five of the experts.
951 This makes it possible to check the consistency of the reviewers, to ensure
952 that they are each reaching comparable decisions regarding whether two (or
953 more) records refer to the same individual or not.

954 After an initial round of clustering, subsets of these combined records were
955 then re-examined to identify previously missed groups of records that refer
956 to the same individual, particularly across years (e.g., records with dates of
957 death 2011/12/31 and 2012/01/01 might refer to the same individual) and
958 governorates (e.g., records with neighboring locations of death might refer
959 to the same individual).

APPENDIX B: UNIQUE ENTITY ESTIMATION PROOFS

First, we introduce four indicators. First, let \mathbb{I}_2 denote every 2-vertex clique in G^* (recall that G^* is the original graph and G' is the observed one):

$$(10) \quad \mathbb{I}_2 = \begin{cases} 1, & \text{if this clique is in } G'. \\ 0, & \text{otherwise.} \end{cases}$$

Second, let \mathbb{I}_{33} denote every 3-vertex clique in G^* :

$$(11) \quad \mathbb{I}_{33} = \begin{cases} 1, & \text{if this clique remains as a 3-clique in } G'. \\ 0, & \text{otherwise.} \end{cases}$$

Third, let \mathbb{I}_{32} denote every 3-vertex clique in G^* :

$$(12) \quad \mathbb{I}_{32} = \begin{cases} 1, & \text{if this clique breaks to a 2-clique in } G'. \\ 0, & \text{otherwise.} \end{cases}$$

Finally, let \mathbb{I}_{31} denote every 3-vertex clique in G^* :

$$(13) \quad \mathbb{I}_{31} = \begin{cases} 1, & \text{if this clique breaks into only 1-cliques in } G'. \\ 0, & \text{otherwise.} \end{cases}$$

B.1. Expectation. We now prove that our estimator is unbiased. Consider

$$(14) \quad \mathbb{E}[n'_3] = \mathbb{E}\left[\sum_{i=1}^{n_3^*} \mathbb{I}_{33i}\right] = n_3^* \cdot p^2 \cdot (3 - 2p),$$

$$(15) \quad \begin{aligned} \mathbb{E}[n'_2] &= \mathbb{E}\left[\sum_{i=1}^{n_2^*} \mathbb{I}_{2i}\right] + \mathbb{E}\left[\sum_{i=1}^{n_3^*} \mathbb{I}_{32i}\right] \\ &= n_2^* \cdot p + n_3^* \cdot (3 \cdot (1 - p)^2 \cdot p), \quad \text{and} \end{aligned}$$

$$(16) \quad \begin{aligned} \mathbb{E}[n'_1] &= n_1^* + \mathbb{E}\left[\sum_{i=1}^{n_2^*} (1 - \mathbb{I}_{2i})\right] + \mathbb{E}\left[\sum_{i=1}^{n_3^*} \mathbb{I}_{32i}\right] + \mathbb{E}\left[\sum_{i=1}^{n_3^*} \mathbb{I}_{31i}\right] \\ &= n_1^* + n_2^* \cdot (2 \cdot (1 - p)) + n_3^* \cdot (3 \cdot (1 - p)^2 \cdot p) \\ &\quad + n_3^* \cdot (3 \cdot (1 - p)^3). \end{aligned}$$

Our estimator is unbiased via equations 16, 15, 14:

$$\begin{aligned}
\mathbb{E}[LSHE] &= \mathbb{E}[n'_1 + n'_2 \cdot \frac{2p-1}{p} \\
&\quad + n'_3 \cdot \frac{1-6 \cdot (1-p)^2 \cdot p}{p^2 \cdot (3-2p)} + \sum_{i=4}^M n_i] \\
&= \mathbb{E}[n'_1] + \frac{2p-1}{p} \cdot \mathbb{E}[n'_2] \\
&\quad + \frac{1-6 \cdot (1-p)^2 \cdot p}{p^2 \cdot (3-2p)} \cdot \mathbb{E}[n'_3] + \mathbb{E}[\sum_{i=4}^M n_i] \\
&= n_1^* + n_2^* + n_3^* + \sum_{i=4}^N n_i^* \\
&= n.
\end{aligned}$$

B.2. Variance. We now turn to deriving the variance of our proposed estimator, showing that

$$\mathbb{V}[LSHE] = n_3^* \cdot \frac{(p-1)^2 \cdot (3p^2 - p + 1)}{p^2 \cdot (3-2p)} + n_2^* \cdot \frac{(1-p)}{p}.$$

970 Consider

$$\begin{aligned}
\mathbb{V}[LSHE] &= \mathbb{V}[\frac{1-6 \cdot (1-p)^2 \cdot p}{p^2 \cdot (3-2p)} \cdot \sum_{i=1}^{n_3^*} \mathbb{I}_{33i} \\
&\quad + \frac{2p-1}{p} \cdot (\sum_{i=1}^{n_2^*} \mathbb{I}_{2i} + \sum_{i=1}^{n_3^*} \mathbb{I}_{32i}) \\
(17) \quad &\quad + \sum_{i=1}^{n_3^*} \mathbb{I}_{31i} + \sum_{i=1}^{n_3^*} \mathbb{I}_{32i} + \sum_{i=1}^{n_2^*} (1 - \mathbb{I}_{2i})] \\
&= \mathbb{V}ar[\frac{1-6 \cdot (1-p)^2 \cdot p}{p^2 \cdot (3-2p)} \cdot \sum_{i=1}^{n_3^*} \mathbb{I}_{33i} + \frac{3p-1}{p} \cdot \sum_{i=1}^{n_3^*} \mathbb{I}_{32i} \\
&\quad + 3 \cdot \sum_{i=1}^{n_2^*} \mathbb{I}_{31i} - \frac{1}{p} \cdot \sum_{i=1}^{n_2^*} \mathbb{I}_{2i}].
\end{aligned}$$

971 Next, we replace $1 - 6 \cdot (1-p)^2 \cdot p$ by a , and by simplifying equation 17,

we find

972

$$\begin{aligned}
 &= \frac{a^2}{(p^2 \cdot (3 - 2p))^2} \cdot \mathbb{V}[\sum_{i=1}^{n_3^*} \mathbb{I}_{33i}] + \frac{(3p - 1)^2}{p^2} \cdot \mathbb{V}[\sum_{i=1}^{n_3^*} \mathbb{I}_{32i}] \\
 (18) \quad &+ 9 \cdot \mathbb{V}[\sum_{i=1}^{n_2^*} \mathbb{I}_{31i}] - \frac{1}{p^2} \cdot \mathbb{V}[\sum_{i=1}^{n_2^*} \mathbb{I}_{2i}] + Cov(\sum_{i=1}^{n_3^*} \mathbb{I}_{33i}, \sum_{i=1}^{n_3^*} \mathbb{I}_{32i}) \\
 &+ Cov(\sum_{i=1}^{n_3^*} \mathbb{I}_{33i}, \sum_{i=1}^{n_3^*} \mathbb{I}_{31i}) + Cov(\sum_{i=1}^{n_3^*} \mathbb{I}_{32i}, \sum_{i=1}^{n_3^*} \mathbb{I}_{31i}).
 \end{aligned}$$

Note that the covariance of $\sum_{i=1}^{n_3^*} \mathbb{I}_{2i}$ with any indicator is zero due to independence. Furthermore, since the indicators are Bernoulli distributed random variables, the variance is easily found. Consider

973

974

975

$$(19) \quad \mathbb{V}ar[\sum_{i=1}^{n_3^*} \mathbb{I}_{33i}] = \frac{a^2 \cdot (1 - p^2 \cdot (3 - 2p))}{p^2 \cdot (3 - 2p)} \cdot n_3^*$$

$$(20) \quad \mathbb{V}ar[\sum_{i=1}^{n_3^*} \mathbb{I}_{32i}] = \frac{3 \cdot (3p - 1)^2 \cdot (1 - p)^2 \cdot (1 - 3p \cdot (1 - p)^2)}{p} \cdot n_3^*$$

$$(21) \quad \mathbb{V}ar[\sum_{i=1}^{n_3^*} \mathbb{I}_{31i}] = 9 \cdot (1 - p)^3 \cdot (1 - (1 - p)^3) \cdot n_3^*$$

$$(22) \quad \mathbb{V}ar[\sum_{i=1}^{n_2^*} \mathbb{I}_{2i}] = \frac{(1 - p)}{p} \cdot n_2^*$$

Using equations 19 – 22, the covariance simplifies to

$$\begin{aligned}
 &Cov(\sum_{i=1}^{n_3^*} \mathbb{I}_{33i}, \sum_{i=1}^{n_3^*} \mathbb{I}_{32i}) \\
 &= \mathbb{E}[\sum_{i=1}^{n_3^*} \mathbb{I}_{33i} \cdot \sum_{i=1}^{n_3^*} \mathbb{I}_{32i}] - \mathbb{E}[\sum_{i=1}^{n_3^*} \mathbb{I}_{33i}] \mathbb{E}[\sum_{i=1}^{n_3^*} \mathbb{I}_{32i}] \\
 &= \sum_{i=1}^{n_3^*} \sum_{j=1}^{n_3^*} \mathbb{E}[\mathbb{I}_{33i} \cdot \mathbb{I}_{32j}] - \mathbb{E}[\mathbb{I}_{33i}] \mathbb{E}[\mathbb{I}_{32j}] \\
 &= -6 \cdot a \cdot n_3^*
 \end{aligned}$$

When $i = j$, since \mathbb{I}_{33j} and \mathbb{I}_{32j} are mutually exclusive, $\mathbb{E}[\mathbb{I}_{33i} \cdot \mathbb{I}_{32j}] = 0$. Otherwise when $i \neq j$, \mathbb{I}_{33j} and \mathbb{I}_{32j} are independent and $\mathbb{E}[\mathbb{I}_{33i} \cdot \mathbb{I}_{32j}] = 0$. Similarly,

$$Cov\left(\sum_{i=1}^{n_3^*} \mathbb{I}_{33i}, \sum_{i=1}^{n_1^*} \mathbb{I}_{32i}\right) = -6 \cdot a \cdot (1-p)^3 \cdot n_3^*$$

and

$$Cov\left(\sum_{i=1}^{n_2^*} \mathbb{I}_{32i}, \sum_{i=1}^{n_1^*} \mathbb{I}_{31i}\right) = -18 \cdot (1-p)^5 \cdot (3p-1) \cdot n_3^*.$$

It then follows that

$$\begin{aligned} \mathbb{V}ar[LSHE] &= n_3^* \cdot \left(\frac{3 \cdot (3p-1)^2 \cdot (1-p)^2 \cdot (1-3p \cdot (1-p)^2)}{p} \right. \\ &+ \frac{(1-6 \cdot (1-p)^2 \cdot p)^2 \cdot (1-p^2 \cdot (3-2p))}{p^2 \cdot (3-2p)} \\ &- (6 \cdot (1-6 \cdot (1-p)^2 \cdot p) \cdot (3p-1) \cdot (1-p)^2) \\ &+ 9 \cdot (1-p)^6 + 3 \cdot (1-p)^3 \Big) \\ &+ n_2^* \cdot \frac{(1-p)}{p} \\ &= n_3^* \cdot \frac{(p-1)^2 \cdot (3p^2 - p + 1)}{p^2 \cdot (3-2p)} + n_2^* \cdot \frac{(1-p)}{p}. \end{aligned}$$

976 **B.3. Variance Monotonicity.** We now prove the monotonicity of the
977 variance of our estimator.

978 **THEOREM 4.** *$\mathbb{V}ar[LSHE]$ is monotonically decreasing when p increases*
979 *in range $(0, 1]$.*

980 **PROOF.**

981 **LEMMA 2.** *First order derivative of $\mathbb{V}ar[LSHE]$ is negative when $p \in$*
982 *$(0, 1]$.*

PROOF. Consider

$$\frac{d(\mathbb{V}[LSHE])}{dp} = \frac{3(2-p)(p-1)(p+1)((p-1)^2 + p^2)}{p^3(2p-3)^2} \cdot n_3^* - p^2 \cdot n_2^*.$$

When $p \in (0, 1]$, $-p^2 < 0$. Because $(2-p)$, $(p+1)$, $(p-1)^2 + p^2$, $p^3(2p-3)^2$ are all positive and $(p-1)$ is the only term that is negative,

$$\frac{3(2-p)(p-1)(p+1)((p-1)^2 + p^2)}{p^3(2p-3)^2} < 0.$$

Thus, $\frac{d(\mathbb{V}[LSHE])}{dp} < 0$. □ 983

By using Lemma 2, we can consequently prove Theorem 4. We also note 984
that when $p = 1$, $\mathbb{V}ar[LSHE] = 0$. □ 985

REFERENCES

- 986 ALEKSANDROV, P. S. (1956). *Combinatorial topology* 1. Courier Corporation.
- 987 ANDONI, A. and INDYK, P. (2004). E2lsh: Exact Euclidean Locality Sensitive Hashing
988 Technical Report.
- 989 BAXTER, R., CHRISTEN, P., CHURCHES, T. et al. (2003). A comparison of fast blocking
990 methods for record linkage. In *ACM SIGKDD* 3 25–27.
- 991 BHATTACHARYA, I. and GETOOR, L. (2006). A Latent Dirichlet Model for Unsupervised
992 Entity Resolution. In *SDM* 5 59. SIAM.
- 993 BRODER, A. Z. (1997a). On the resemblance and containment of documents. In *Compression
994 and Complexity of Sequences 1997. Proceedings* 21–29. IEEE.
- 995 BRODER, A. Z. (1997b). On the Resemblance and Containment of Documents. In *the
996 Compression and Complexity of Sequences* 21–29.
- 997 CHAZELLE, B., RUBINFELD, R. and TREVISAN, L. (2005). Approximating the minimum
998 spanning tree weight in sublinear time. *SIAM Journal on computing* 34 1370–1379.
- 999 CHRISTEN, P. (2012). A survey of indexing techniques for scalable record linkage and
1000 deduplication. *IEEE Transactions on Knowledge and Data Engineering* 24 1537–1555.
- 1001 CHRISTEN, P. (2014). Preparation of a real voter data set for record linkage and duplicate
1002 detection research.
- 1003 DEMING, W. E. and GLASSER, G. J. (1959). On the problem of matching lists by samples.
1004 *Journal of the American Statistical Association* 54 403–415.
- 1005 ERDOS, P. and RÉNYI, A. (1960). On the evolution of random graphs. *Publ. Math. Inst.
1006 Hung. Acad. Sci* 5 17–60.
- 1007 FELLEGI, I. and SUNTER, A. (1969). A Theory for Record Linkage. *Journal of the Amer-
1008 ican Statistical Association* 64 1183–1210.
- 1009 FRANK, O. (1978). Estimation of the Number of Connected Components in a Graph by
1010 Using a Sampled Subgraph. *Scandinavian Journal of Statistics* 5 177–188.
- 1011 GIONIS, A., INDYK, P., MOTWANI, R. et al. (1999). Similarity search in high dimensions
1012 via hashing. In *Very Large Data Bases (VLDB)* 99 518–529.
- 1013 GRILLO, C. (2016). Judges in Habre Trial Cite HRDAG Analysis.
- 1014 GUTMAN, R., AFENDULIS, C. C. and ZASLAVSKY, A. M. (2013). A Bayesian procedure
1015 for file linking to analyze end-of-life medical costs. *Journal of the American Statistical
1016 Association* 108 34–47.
- 1017 INDYK, P. and MOTWANI, R. (1998). Approximate Nearest Neighbors: Towards Removing
1018 the Curse of Dimensionality. In *STOC* 604–613.
- 1019 LIANG, H., WANG, Y., CHRISTEN, P. and GAYLER, R. (2014). Noise-tolerant approxi-
1020 mate blocking for dynamic real-time entity resolution. In *Pacific-Asia Conference on
1021 Knowledge Discovery and Data Mining* 449–460. Springer.
- 1022 LISEO, B. and TANCREDI, A. (2013). Some advances on Bayesian record linkage and infer-
1023 ence for linked data. URL [http://www.ine.es/e/essnetdi.ws2011/ppts/Liseo.Tancredi.
1024 pdf](http://www.ine.es/e/essnetdi.ws2011/ppts/Liseo.Tancredi.pdf).
- 1025 LUO, C. and SHRIVASTAVA, A. (2017). Arrays of (locality-sensitive) Count Estimators
1026 (ACE): High-Speed Anomaly Detection via Cache Lookups. *CoRR* abs/1706.06664.
- 1027 MCCALLUM, A., NIGAM, K. and UNGAR, L. H. (2000). Efficient clustering of high-
1028 dimensional data sets with application to reference matching. In *Proceedings of the
1029 sixth ACM SIGKDD international conference on Knowledge discovery and data mining*
1030 169–178. ACM.
- 1031 MCCALLUM, A. and WELLNER, B. (2004). Conditional Models of Identity Uncertainty
1032 with Application to Noun Coreference. In *Advances in Neural Information Processing
1033 Systems (NIPS '04)* 905–912. MIT Press.

- PAULEVÉ, L., JÉGOU, H. and AMSALEG, L. (2010). Locality sensitive hashing: A comparison of hash function types and querying mechanisms. *Pattern Recognition Letters* **31** 1348–1358.
- PRICE, M. and BALL, P. (2015a). Selection bias and the statistical patterns of mortality in conflict. *Statistical Journal of the IAOS* **31** 263–272.
- PRICE, M. and BALL, P. (2015b). The Limits of Observation for Understanding Mass Violence. *Canadian Journal of Law and Society/Revue Canadienne Droit et Société* **30** 237–257.
- PRICE, M., GOHDES, A. and BALL, P. (2015). Documents of war: Understanding the Syrian conflict. *Significance* **12** 14–19.
- PRICE, M., KLINGNER, J., QTIESH, A. and BALL, P. (2013). Updated statistical analysis of documentation of killings in the Syrian Arab Republic. *United Nations Office of the UN High Commissioner for Human Rights*.
- PRICE, M., KLINGNER, J., QTIESH, A. and BALL, P. (2014). Updated statistical analysis of documentation of killings in the Syrian Arab Republic. *United Nations Office of the UN High Commissioner for Human Rights*.
- PROVAN, J. S. and BALL, M. O. (1983). The Complexity of Counting Cuts and of Computing the Probability that a Graph is Connected. *SIAM Journal on Computing* **12** 777–788.
- RAJARAMAN, A. and ULLMAN, J. D. (2012). *Mining of massive datasets*. Cambridge University Press.
- SADINLE, M. et al. (2014). Detecting duplicates in a homicide registry using a Bayesian partitioning approach. *The Annals of Applied Statistics* **8** 2404–2434.
- SADOSKY, P., SHRIVASTAVA, A., PRICE, M. and STEORTS, R. C. (2015). Blocking Methods Applied to Casualty Records from the Syrian Conflict. *arXiv preprint arXiv:1510.07714*.
- SHRIVASTAVA, A. (2017). Optimal Densification for Fast and Accurate Minwise Hashing. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- SHRIVASTAVA, A. and LI, P. (2014a). Densifying one permutation hashing via rotation for fast near neighbor search. In *Proceedings of The 31st International Conference on Machine Learning* 557–565.
- SHRIVASTAVA, A. and LI, P. (2014b). Improved Densification of One Permutation Hashing. In *Proceedings of The 30th Conference on Uncertainty in Artificial Intelligence*.
- SHRIVASTAVA, A. and LI, P. (2014c). In Defense of Minhash over Simhash. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics* 886–894.
- SPRING, R. and SHRIVASTAVA, A. (2017a). Scalable and sustainable deep learning via randomized hashing. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* 445–454. ACM.
- SPRING, R. and SHRIVASTAVA, A. (2017b). A New Unbiased and Efficient Class of LSH-Based Samplers and Estimators for Partition Function Computation in Log-Linear Models. *arXiv preprint arXiv:1703.05160*.
- STEORTS, R. C., HALL, R. and FIENBERG, S. E. (2016). A Bayesian Approach to Graphical record Linkage and De-duplication. *Journal of the American Statistical Society*.
- STEORTS, R. C., VENTURA, S. L., SADINLE, M. and FIENBERG, S. E. (2014). A Comparison of Blocking Methods for Record Linkage. In *International Conference on Privacy in Statistical Databases* 253–268.
- TANCREDI, A. and LISEO, B. (2011). A hierarchical Bayesian approach to record linkage and population size problems. *Annals of Applied Statistics* **5** 1553–1585.
- VATSALAN, D., CHRISTEN, P., O’KEEFE, C. M. and VERYKIOS, V. S. (2014). An eval-

- 1084 uation framework for privacy-preserving record linkage. *Journal of Privacy and Confi-*
1085 *dentiality* **6** 3.
- 1086 WANG, Y., SHRIVASTAVA, A. and RYU, J. (2017). FLASH: Randomized Algorithms Accel-
1087 erated over CPU-GPU for Ultra-High Dimensional Similarity Search. *ArXiv e-prints*.
- 1088 WINKLER, W. E. (2004). Approximate String Comparator Search Strategies for Very
1089 Large Administrative Lists. *Proceedings of the Section on Survey Research Methods,*
1090 *American Statistical Association*.
- 1091 WINKLER, W. E. (2006). Overview of record linkage and current research directions. In
1092 *Bureau of the Census*. Citeseer.

1093 DEPARTMENT OF COMPUTER SCIENCE
RICE UNIVERSITY, HOUSTON, TX USA
E-MAIL: beidi.chen@rice.edu; anshumali@rice.edu

DEPARTMENT OF STATISTICAL SCIENCE
AND COMPUTER SCIENCE
DUKE UNIVERSITY
DURHAM, NC USA
E-MAIL: beka@stat.duke.edu