

# Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)

Brian Kunding<sup>\*</sup>, Jerome P. Reiter<sup>\*</sup> and Rebecca C. Steorts<sup>\*</sup>

**Abstract.** Within the field of record linkage, Bayesian methods have the crucial advantage of quantifying uncertainty from imperfect linkages. However, current implementations of Bayesian Fellegi-Sunter models are computationally intensive, making them challenging to use on larger-scale record linkage tasks. To address these computational difficulties, we propose fast beta linkage (**fabl**), an extension to the Beta Record Linkage (**BRL**) method of [Sadinle \(2017\)](#). Specifically, we use independent prior distributions over the matching space, allowing us to use hashing techniques that reduce computational overhead. This also allows us to complete pairwise record comparisons over large data files through parallel computing and to reduce memory costs through a new technique called storage efficient indexing. Through simulations and two case studies, we show that **fabl** can have markedly increased speed with minimal loss of accuracy when compared to **BRL**.

**Keywords:** data fusion, data cleaning, entity resolution, hashing, record linkage.

## 1 Introduction

Before conducting data analysis, it is often necessary to identify duplicate records across two data files. This is an increasingly important task in “data cleaning” and is used for inferential and predictive analyses in fields such as statistics, computer science, machine learning, political science, economics, precision medicine, official statistics, and others (e.g., [Christen, 2012](#); [Gutman et al., 2013](#); [Dalzell and Reiter, 2018](#); [Tang et al., 2020](#)). In this article, we consider bipartite record linkage, which merges two data files that contain duplications across, but not within, the respective data files.

Many probabilistic record linkage methods rely on the seminal work of [Fellegi and Sunter \(1969\)](#) and [Newcombe et al. \(1959\)](#). In their approach, the data analyst first creates comparison vectors for each pair of records in the data files. These vectors indicate how similar the records are on a set of variables measured in both files, known as the linkage variables. Using these comparison vectors, the analyst classifies each pair as a match or nonmatch using a likelihood ratio test. Crucially, these decisions are made independently for each pair. The [Fellegi and Sunter \(1969\)](#) approach has been extended for a wide variety of applications (e.g., [Winkler and Thibaudeau, 1990](#); [Fair, 2004](#); [Wagner et al., 2014](#); [Gill and Goldacre, 2003](#); [Enamorado et al., 2019](#); [Aleshin-Guendel and Sadinle, 2023](#)). An alternative paradigm is to model the linkage variables directly (e.g., [Tancredi et al., 2011](#); [Steorts et al., 2016](#); [Marchant et al., 2021](#); [Betancourt et al., 2022](#)). In this article, we build on the contributions to the comparison vector approach.

---

<sup>\*</sup>Department of Statistical Science, Duke University, P.O. Box 90251, Durham, NC 27708, USA  
[brian.kunding@duke.edu](mailto:brian.kunding@duke.edu), [jreiter@duke.edu](mailto:jreiter@duke.edu), [beka@stat.duke.edu](mailto:beka@stat.duke.edu)

## 2 *Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)*

37 The independent pairwise matching assumption from [Fellegi and Sunter \(1969\)](#) is  
38 popular for its mathematical simplicity. However, in many situations, there are no  
39 duplications within a data file, meaning that each record in one file should be linked  
40 with at most one other record in the other file. Thus, when the procedure results in  
41 many-to-one matches, some of these links must be false. Analysts typically use an  
42 additional post-processing step to turn the list of linked records into a bipartite matching  
43 ([Jaro, 1989](#)), but this model misspecification can still lead to poor results ([Sadinle, 2017](#)).

44 Alternatively, analysts can embed one-to-one matching requirements into the model  
45 specification, at an additional computational cost. [Larsen \(2005\)](#) employed a Metropolis-  
46 Hastings algorithm to only allow sampling matches that respect one-to-one assumptions,  
47 but such algorithms exhibit slow mixing due to the combinatorial nature of the con-  
48 strained matching space. [Fortunato \(2010\)](#) used simulated annealing to target the space  
49 of matches permitted under the one-to-one constraint, but the method is computationally  
50 intensive and, to our knowledge, has not been applied on data files with more than  
51 100 records. [Sadinle \(2017\)](#) proposed the Beta Record Linkage model (BRL), using a  
52 prior distribution over the space of bipartite matchings to strictly enforce one-to-one  
53 requirements throughout a Gibbs sampler. BRL has been shown to work on larger tasks  
54 than previous one-to-one methods, but in our experience, it becomes slow when applied  
55 to files with more than a few thousand records.

56 In this article, we propose fast beta linkage (**fabl**), which extends the BRL model  
57 for increased efficiency and scalability. We use independent prior distributions for the  
58 matching statuses of each record, and modify the decision theoretic technique of [Sadinle](#)  
59 [\(2017\)](#) to ensure our linkage estimates are bipartite. This approach allows us to (1)  
60 employ hashing techniques that speed up calculations and reduce computational costs,  
61 (2) compute the pairwise record comparisons over large data files via parallel computing,  
62 and (3) reduce memory costs through what we call storage efficient indexing. These  
63 contributions allow **fabl** to perform record linkage on much larger data files than  
64 previous Bayesian Fellegi-Sunter models at significantly increased speed with similar  
65 levels of accuracy. In particular, computation time under BRL grows quadratically, with  
66 the size of each data file, while computation time under **fabl** grows linearly, only with  
67 the size of the smaller data file.

68 In what follows, Section 2 reviews the work of [Fellegi and Sunter \(1969\)](#) and [Sadinle](#)  
69 [\(2017\)](#). Section 3 proposes the **fabl** model, provides the Gibbs sampler for posterior  
70 inference, and describes the loss function used to calculate the Bayes estimate for the  
71 bipartite matching. Section 4 introduces the hashing technique and storage efficient  
72 indexing used to increase the speed of calculations and the scale of linkage tasks amenable  
73 to **fabl**. Sections 5 and 6 demonstrate the speed and accuracy of **fabl** through simulation  
74 studies and case studies of homicides from the El Salvadoran Civil War and the National  
75 Long Term Care Study. Finally, Section 7 summarizes our contributions and highlights  
76 areas for further research.

## 2 Review of Fellegi-Sunter Approaches for Record Linkage

Consider two data files  $A$  and  $B$  comprising  $n_A$  and  $n_B$  records, respectively, and including  $F$  linkage variables measured in both files. For  $i = 1, \dots, n_A$ , let record  $i$  be given by  $A_i = (A_{i1}, \dots, A_{iF})$ , so that  $A = (A_i : i = 1, \dots, n_A)$ . Similarly, for  $j = 1, \dots, n_B$ , let record  $j$  be given by  $B_j = (B_{j1}, \dots, B_{jF})$ , so that  $B = (B_j : j = 1, \dots, n_B)$ . Without loss of generality, denote files such that  $n_A \geq n_B$ .

Intuitively, matching records (those that refer to the same entity) should have similar values of the linking variables; records that are nonmatching should have dissimilar values. Fellegi and Sunter (1969) proposed encoding this using a comparison vector  $\gamma_{ij}$  computed for each record pair  $(i, j)$  in  $A \times B$ . Specifically, they define  $\gamma_{ij} = (\gamma_{ij}^1, \dots, \gamma_{ij}^F)$ , where each  $\gamma_{ij}^f$  is a value indicating the similarity of field  $f$  for records  $A_i$  and  $B_j$ . We define the comparison matrix  $\gamma \in \mathbb{R}^{n_A n_B \times F}$  as the collection of all record pairs  $\gamma_{ij}$ .

When linking variable  $f$  is categorical, a common way to define  $\gamma_{ij}^f$  is an indicator for exact agreement. For example, if zip code is linking variable  $f$ , we can set  $\gamma_{ij}^f = 1$  when the zip codes for records  $A_i$  and  $B_j$  agree exactly, and set  $\gamma_{ij}^f = 2$  when they do not. For numerical linking variables, we can use the absolute difference of the two values. For example, if age is linking variable  $f$ , we can set  $\gamma_{ij}^f = 1$  when the ages for records  $A_i$  and  $B_j$  match exactly,  $\gamma_{ij}^f = 2$  when the ages are within one year but not equal, and  $\gamma_{ij}^f = 3$  when the ages are two or more years apart. For text variables like names, we can use string distance metrics such as Levenstein or Jaro-Winkler distance (Cohen et al., 2003). We then set thresholds that allow us to represent comparisons through discrete levels of disagreement (Bilenko and Mooney, 2006; Elmagarmid et al., 2007).

More generally, let  $\mathcal{S}_f(i, j)$  denote a similarity measure for linking variable  $f$  of records  $A_i$  and  $B_j$ . The range of  $\mathcal{S}_f$  can be divided into  $L_f$  intervals denoted by  $I_{f1}, \dots, I_{fL_f}$ . Here,  $I_{f1}$  represents the highest level of agreement (including complete agreement) and  $I_{fL_f}$  represents the highest level of disagreement (including complete disagreement). Thus, we can construct comparison vectors such that  $\gamma_{ij}^f = l$  if  $\mathcal{S}_f(i, j) \in I_{fl}$ . The choices of  $I_{fl}$  are application specific, as we discuss in the simulation and case studies.

In the construction of comparison vectors, it is common to encounter missing information in record  $A_i$  or  $B_j$ . As a result, the comparison vector  $\gamma_{ij}$  will have missing values. We assume that this missingness occurs completely at random (MCAR, per Little and Rubin (2002)). To notate a missing value in any  $\gamma_{ij}^f$ , we use  $I_{obs}(\gamma_{ij}^f) = 1$  when  $\gamma_{ij}^f$  is observed and  $I_{obs}(\gamma_{ij}^f) = 0$  otherwise. With the MCAR assumption, we can marginalize over the missing data, and do all computation simply using the observed data.

Having defined comparison vectors, we now turn to the Fellegi and Sunter (1969) model for record linkage. We begin with notation for defining linkages. Under bipartite matching, the set of matches across  $A$  and  $B$  can be represented in two equivalent ways.

#### 4 Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)

First, we may use a matrix  $\Delta \in \{0, 1\}^{n_A \times n_B}$ , where

$$\Delta_{ij} = \begin{cases} 1, & \text{if records } A_i \text{ and } B_j \text{ refer to the same entity;} \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

This sparse matrix representation can become cumbersome for large linkage tasks. More compactly, bipartite matching can also be viewed as a labeling  $\mathbf{Z} = (Z_1, \dots, Z_{n_B})$  for the records in  $B$  such that

$$Z_j = \begin{cases} i, & \text{if records } A_i \text{ and } B_j \text{ refer to the same entity;} \\ n_A + j, & \text{if record } B_j \text{ does not have a match in } A. \end{cases} \quad (2)$$

112 We can go back and forth between the two using  $\Delta_{ij} = I(Z_j = i)$ , where  $I(\cdot) = 1$  when  
 113 the expression inside the parentheses is true, and  $I(\cdot) = 0$  otherwise. When presenting  
 114 the models, we use both representations for convenience.

Let  $\Gamma_{ij}$  represent a random variable for the comparison vector for arbitrary  $A_i$  and  $B_j$ . Thus,  $\gamma_{ij}$  is a realization of  $\Gamma_{ij}$ . For modeling the collection of  $n_A n_B$  random variables  $\Gamma_{ij}$ , Fellegi and Sunter (1969) employ two independence assumptions: first, that comparison vectors are independent given the matching status of the record pair, and second, that the matching status of each record pair is independent of the matching status of other pairs. Using these independence assumptions, one specifies a mixture model for each  $\Gamma_{ij}$  (e.g., as in Winkler, 1999; Jaro, 1989; Larsen and Rubin, 2001; Enamorado et al., 2019). We have

$$\Gamma_{ij} \mid \Delta_{ij} = 1 \stackrel{iid}{\sim} \mathcal{M}(\mathbf{m}), \quad (3a)$$

$$\Gamma_{ij} \mid \Delta_{ij} = 0 \stackrel{iid}{\sim} \mathcal{U}(\mathbf{u}), \quad (3b)$$

$$\Delta_{ij} \stackrel{iid}{\sim} \text{Bernoulli}(\lambda). \quad (3c)$$

115 Here,  $\mathcal{M}$  and  $\mathcal{U}$  are the distributions for matching and nonmatching record pairs,  $\mathbf{m}$   
 116 and  $\mathbf{u}$  are their respective sets of parameters, and  $\lambda$  is the marginal probability that  
 117 a record pair is a match. When using comparison vectors with discrete agreement  
 118 levels,  $\mathcal{M}$  and  $\mathcal{U}$  are collections of independent multinomial distributions for each  
 119 linkage feature. Accordingly,  $\mathbf{m} = (\mathbf{m}_1, \dots, \mathbf{m}_F)$ , where  $\mathbf{m}_f = (m_{f1}, \dots, m_{fL_f})$  and  
 120  $m_{fl} = p(\Gamma_{ij}^f = l \mid \Delta_{ij} = 1)$  for all fields  $f$  and agreement levels  $l$ . The  $\mathbf{u}$  parameters are  
 121 defined similarly, with  $u_{fl} = p(\Gamma_{ij}^f = l \mid \Delta_{ij} = 0)$ .

Sadinle (2017) presents a Bayesian version of the model in (3a) through (3c). He uses uniform Dirichlet prior distributions for each  $\mathbf{m}_f$  and  $\mathbf{u}_f$ , and replaces (3c) with a prior distribution for  $\mathbf{Z}$  that he calls the “beta distribution for bipartite matching.” This assigns a Bernoulli distribution for the indicator that a record in  $B$  has a match in  $A$ , that is,  $I(Z_j \leq n_A) \sim \text{Bernoulli}(\pi)$ . Additionally,  $\pi \sim \text{Beta}(\alpha_\pi, \beta_\pi)$ , where  $\alpha_\pi$  and  $\beta_\pi$  are known hyperparameters. It follows that the number of records in  $B$  that have matches, denoted  $n_{AB}(\mathbf{Z}) = \sum_{j=1}^{n_B} I(Z_j \leq n_A)$ , is distributed according to a Beta-Binomial( $n_B, \alpha_\pi, \beta_\pi$ ). Conditioning on the set of records in  $B$  that have matches,

formally denoted  $\{I(Z_j \leq n_A)\}_{j=1}^{n_B}$ , all  $n_A!/(n_A - n_{AB}(\mathbf{Z}))!$  bipartite matchings are taken to be equally likely. Thus, the prior distribution used by [Sadinle \(2017\)](#) is given by

$$p(\mathbf{Z}|\alpha_\pi, \beta_\pi) = \frac{(n_A - n_{AB}(\mathbf{Z}))!}{n_A!} \frac{B(n_{AB}(\mathbf{Z}) + \alpha_\pi, n_B - n_{AB}(\mathbf{Z}) + \beta_\pi)}{B(\alpha_\pi, \beta_\pi)}, \quad (4)$$

where  $B(\cdot, \cdot)$  represents the Beta function. This prior strictly enforces one-to-one matching, inducing a Gibbs sampler that removes previously matched records from the set of candidate records when sampling each  $Z_j$ . This makes the sampler inherently sequential, which can be slow when working on linkage tasks with more than a few thousand records.

### 3 Fast Beta Linkage

Instead of the prior over  $\mathbf{Z}$  from [Sadinle \(2017\)](#), we follow [Wortman \(2019\)](#) and use independent priors for each component  $Z_j$ . However, unlike [Wortman \(2019\)](#) who proposes a flat prior for  $Z_j$ , we use the fast beta linkage (`fabl`) prior below. For each  $Z_j$ , we have

$$p(Z_j = q|\pi) = \begin{cases} \frac{1}{n_A}\pi, & q \leq n_A; \\ 1 - \pi, & q = n_A + j; \end{cases} \quad (5)$$

$$\pi \sim \text{Beta}(\alpha_\pi, \beta_\pi).$$

We can interpret (5) as follows: record  $B_j$  has some match in  $A$  with probability  $\pi$ , and each record  $A_i$  is equally likely to be that match. The hyperparameters  $\alpha_\pi$  and  $\beta_\pi$  encode prior beliefs about the proportion of records in  $B$  that have matches in  $A$ .

In the [Wortman \(2019\)](#) flat prior, each value  $\{1, \dots, n_A, n_A + j\}$  is a priori equally likely for  $Z_j$ . This amounts to a prior probability of  $n_A/(n_A + 1)$  that record  $B_j$  has a match in  $A$ . In our preliminary studies, we have found that this highly informative prior weighting on matching can result in overly high match rates. Hence, we prefer (5). We also note that the flat prior is equivalent to a special case of the fast beta prior with  $\pi$  fixed at the mean of a Beta  $(1, 1/n_A)$  random variable.

Linkage with `fabl` is conducted at the record level, rather than at the record pair level, as in the [Fellegi and Sunter \(1969\)](#) model. That is,  $\pi$  in (5) under `fabl` estimates the proportion of records in  $B$  that have matches, whereas  $\lambda$  in (3c) in the [Fellegi and Sunter \(1969\)](#) model estimates the proportion of record pairs that are matches. In the bipartite case, we conjecture that some analysts will find  $\pi$  to be a more interpretable parameter than  $\lambda$ . In this setting, there are at most  $n_B$  matching pairs out of  $n_A n_B$  total pairs, meaning that  $\lambda$  is bounded above by  $1/n_A$  and tends towards 0 as the size of the linkage task grows. Additionally, while the [Fellegi and Sunter \(1969\)](#) model makes  $n_A n_B$  independent matching decisions and BRL makes  $n_B$  dependent matching decisions, `fabl` strikes a middle ground between the two, making  $n_B$  independent matching decisions. As shown in Sections 5 and 6, this allows `fabl` to fit a Bayesian record linkage model like BRL while making computational efficiency gains possible by exploiting independence.

## 6 Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)

For clarity, we present the full model for **fabl** below:

$$\mathcal{L}(\mathbf{Z}, \mathbf{m}, \mathbf{u} \mid \gamma) = \prod_{i=1}^{n_A} \prod_{j=1}^{n_B} \prod_{f=1}^F \prod_{l=1}^{L_f} \left[ m_{fl}^{I(Z_j=i)} u_{fl}^{I(Z_j \neq i)} \right]^{I(\gamma_{ij}^f=l)I_{obs}(\gamma_{ij}^f)}, \quad (6a)$$

$$\mathbf{m}_f \sim \text{Dirichlet}(\alpha_{f1}, \dots, \alpha_{fL_f}), \forall f = 1, \dots, F, \quad (6b)$$

$$\mathbf{u}_f \sim \text{Dirichlet}(\beta_{f1}, \dots, \beta_{fL_f}), \forall f = 1, \dots, F, \quad (6c)$$

$$p(Z_j = q \mid \pi) = \begin{cases} \frac{1}{n_A} \pi, & q \leq n_A; \\ 1 - \pi, & q = n_A + j; \end{cases} \quad (6d)$$

$$\pi \sim \text{Beta}(\alpha_\pi, \beta_\pi). \quad (6e)$$

We estimate the posterior distribution of the parameters in (6a - 6e) using a Gibbs sampler. Supplement A presents derivations for the full conditional distributions. We initialize  $\mathbf{m}$ ,  $\mathbf{u}$  and  $\pi$  from random draws from their prior distributions, and initialize  $\mathbf{Z}$  to reflect no matches across data files; that is,  $\mathbf{Z} = (n_A + 1, \dots, n_A + n_B)$ . Denote the current draw of the sampler by  $(s)$  and the updated draw by  $(s+1)$ . To update  $\mathbf{m}_f$  and  $\mathbf{u}_f$  for all  $f = 1, \dots, F$ , we sample

$$\mathbf{m}_f^{(s+1)} \mid \gamma, \mathbf{Z}^{(s)}, \mathbf{u}^{(s)}, \pi^{(s)} \sim \text{Dirichlet}(\alpha_{f1}(\mathbf{Z}^{(s)}), \dots, \alpha_{fL_f}(\mathbf{Z}^{(s)})), \quad (7a)$$

$$\mathbf{u}_f^{(s+1)} \mid \gamma, \mathbf{Z}^{(s)}, \mathbf{m}^{(s)}, \pi^{(s)} \sim \text{Dirichlet}(\beta_{f1}(\mathbf{Z}^{(s)}), \dots, \beta_{fL_f}(\mathbf{Z}^{(s)})), \quad (7b)$$

$$\text{where } \alpha_{fl}(\mathbf{Z}^{(s)}) = \alpha_{fl} + \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} I_{obs}(\gamma_{ij}^f) I(\gamma_{ij}^f = l) I(Z_j^{(s)} = i), \quad (7c)$$

$$\text{and } \beta_{fl}(\mathbf{Z}^{(s)}) = \beta_{fl} + \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} I_{obs}(\gamma_{ij}^f) I(\gamma_{ij}^f = l) I(Z_j^{(s)} \neq i). \quad (7d)$$

Next, we sample  $\pi$  from its full conditional, given by

$$\pi^{(s+1)} \mid \gamma, \mathbf{Z}^{(s)}, \mathbf{m}^{(s+1)}, \mathbf{u}^{(s+1)} \sim \text{Beta}(n_{AB}(\mathbf{Z}^{(s)}) + \alpha_\pi, n_B - n_{AB}(\mathbf{Z}^{(s)}) + \beta_\pi), \quad (8)$$

152 where  $n_{AB}(\mathbf{Z}^{(s)}) = \sum_{j=1}^{n_B} I(Z_j^{(s)} \leq n_A)$ .

Lastly, we sample  $\mathbf{Z}$  componentwise from the full conditional for each  $Z_j$ :

$$p\left(Z_j^{(s+1)} = q \mid \gamma, \mathbf{m}^{(s+1)}, \mathbf{u}^{(s+1)}, \pi^{(s+1)}\right) \propto \begin{cases} \frac{\pi^{(s+1)}}{n_A} w_{qj}^{(s+1)}, & q \leq n_A; \\ 1 - \pi^{(s+1)}, & q = n_A + j, \end{cases} \quad (9)$$

where, for all  $i \in \{1, \dots, n_A\}$  and  $j \in \{1, \dots, n_B\}$ ,

$$w_{ij}^{(s)} = \prod_{f=1}^F \prod_{l=1}^{L_f} \left( \frac{m_{fl}^{(s)}}{u_{fl}^{(s)}} \right)^{I(\gamma_{ij}^f=l)I_{obs}(\gamma_{ij}^f)}. \quad (10)$$

153 Finally, to obtain an estimate  $\hat{\mathbf{Z}}$  of the linkage structure, we use the loss functions  
154 and Bayes estimate from [Sadinle \(2017\)](#). Since (5) does not strictly enforce one-to-one

155 matching, it is possible for this Bayes estimate to link multiple records in  $B$  to one record  
 156 in  $A$ . To obtain a Bayes estimate that fulfills the bipartite requirement, we minimize the  
 157 expected loss subject to the constraint that  $\hat{Z}_j \neq \hat{Z}_{j'}$  for all  $j \neq j'$ . See Supplement B  
 158 for details regarding the initial Bayes estimate and this post-processing procedure.

## 159 4 Efficient and Scalable Implementation

160 The scale of linkage tasks possible through BRL is limited by the memory costs of storing  
 161  $n_A n_B$  comparison vectors for every pair of records across the two data files, and the  
 162 speed of the linkage algorithm over those comparison vectors. One approach to reduce  
 163 the number of comparisons is blocking, which places similar records into partitions, or  
 164 “blocks” (Christen, 2019). In deterministic blocking, the modeler chooses fields thought  
 165 to be highly reliable and only compares records that agree on those fields. The record  
 166 linkage method is applied independently across all blocks, which can be done in parallel  
 167 for additional speed gains. Of note, blocking on an unreliable field can lead to missed  
 168 matches, making this form of blocking often undesirable (Steorts et al., 2014).

169 After computing all comparison vectors within a block, the modeler can further  
 170 reduce the number of comparison vectors used in the linkage algorithm through indexing.  
 171 For example, one might only consider pairs with a certain similarity score on a field  
 172 deemed to be important, like last name, or pairs that exactly match on a specified  
 173 number of fields. However, the impact of indexing on model parameters is not well  
 174 understood; Murray (2016) reviewed this issue in the context of the classical Fellegi and  
 175 Sunter (1969) model, leaving the effect of indexing on Bayesian record linkage models to  
 176 future work.

177 With **fabl**, we introduce two techniques to further expand the scalability of prob-  
 178 abilistic record linkage. First, we propose hashing methods that allow us to compute  
 179 summary statistics that reduce the computational complexity of the Gibbs sampler  
 180 and the memory requirements of storing the comparison vectors. Second, we introduce  
 181 storage efficient indexing, which reduces the memory costs associated with unlikely  
 182 matches.

### 183 4.1 Data Representation, Hashing, and Storage

184 Since each component  $\gamma_{ij}^f$  is discrete, there are only finitely many possible realizations of  
 185 the comparison vector  $\gamma_{ij}$ . Let  $P$  be the number of patterns realized in  $\gamma$ . It is bounded  
 186 above by  $P^* = \prod_{f=1}^F (L_f + 1)$ , where the addition of 1 to  $L_f$  for each field accounts for  
 187 the possibility of missing values. This quantity is determined by  $F$  and  $L_f$ , and does  
 188 not scale with  $n_A$  or  $n_B$ .

To obtain a memory efficient representation, we map the agreement pattern of each  
 record pair to a unique integer. Enamorado et al. (2019) accomplished this through a  
 hashing function, which we modify to explicitly handle missing values:

$$h^*(\gamma_{ij}) = \sum_{f=1}^F I_{obs}(\gamma_{ij}^f) 2^{\gamma_{ij}^f + I(f>1) \sum_{d=1}^{f-1} (L_d)}. \quad (11)$$

## 8 Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)

We then map the integers in (11) to sequential integers  $\{1, \dots, P\}$ . Denote each unique agreement pattern as  $h_p = (h_p^1, \dots, h_p^F)$ , and the set of unique agreement patterns as  $\mathcal{P} = \{h_1, \dots, h_P\}$ . When the  $(i, j)$  record pair exhibits agreement pattern  $p$ , we say  $\gamma_{ij} = h_p$ . In calculations, we will at times use the one-hot encoding of agreement pattern  $h_p$ , denoted  $e(h_p)$ . This is a vector of length  $\sum_{f=1}^F L_f$  in which the  $l + \sum_{k=1}^{f-1} L_k$  component is 1 when  $\gamma_{ij}^f = l$ , and 0 otherwise.

For example, consider five fields with binary agreements and possible missingness, denoted  $NA$ . The number of possible patterns is bounded above by  $P^* = 3^5 = 243$ . Suppose that records  $A_5$  and  $B_7$  exhibit agreement pattern  $\gamma_{5,7} = (1, 1, 1, NA, 2)$ , indicating exact agreement on the first three fields, missing information in the fourth field, and disagreement in the fifth field. The expression in (11) gives  $h^*(\gamma_{5,7}) = 2^1 + 2^3 + 2^5 + 0 + 2^{10} = 1066$ . All records with a hashed value of 1066 map to the integer 42. Thus,  $\gamma_{5,7} = h_{42}$ , and this agreement pattern has the one hot encoding  $e(h_{42}) = (1, 0, 1, 0, 1, 0, 0, 0, 0, 1)$ .

We then identify the records in  $A$  with comparison vectors corresponding to each pattern  $p$  for each record  $B_j$ . We denote this set  $r_{p_j} = \{i \in 1, \dots, n_A | \gamma_{ij} = h_p\}$ , and collect all such sets in the nested list  $\mathcal{R} = \{r_{p_j} | p \in \{1, \dots, P\}, j \in \{1, \dots, n_B\}\}$ . We compute the number of records in  $A$  that share agreement pattern  $p$  with record  $B_j$ , given by

$$N_{p_j} = |r_{p_j}| = \sum_{i=1}^{n_A} I(\gamma_{ij} = h_p). \quad (12)$$

We collect these counts in  $\mathcal{N} = \{N_{p_j} | p \in 1, \dots, P, j \in 1, \dots, n_B\}$ .

The set  $\tilde{\gamma} = \{\mathcal{P}, \mathcal{R}, \mathcal{N}\}$  fully characterizes the comparison matrix  $\gamma$  for writing the likelihood function for **fabl**. To see this, we employ the condensed notation

$$m_p = p(\Gamma_{ij} = h_p | Z_j = i) = \prod_{f=1}^F \prod_{l=1}^{L_f} m_{fl}^{I(h_p^f=l)I_{obs}(h_p^f)}, \quad (13a)$$

$$u_p = p(\Gamma_{ij} = h_p | Z_j \neq i) = \prod_{f=1}^F \prod_{l=1}^{L_f} u_{fl}^{I(h_p^f=l)I_{obs}(h_p^f)} \quad (13b)$$

to express the probability that records  $A_i$  and  $B_j$  form agreement pattern  $p$  given that they are a match in (13a), and not a match in (13b). Viewed through the perspective of agreement patterns, the likelihood in (6a) is equivalent to

$$\mathcal{L}(\mathbf{Z}, \mathbf{m}, \mathbf{u} | \tilde{\gamma}) = \prod_{j=1}^{n_B} \prod_{p=1}^P \prod_{i \in r_{p_j}} m_p^{I(Z_j=i)} u_p^{1-I(Z_j=i)}. \quad (14)$$

The likelihood in (14) allows for more efficient posterior inference for  $\mathbf{m}$ ,  $\mathbf{u}$ , and  $\mathbf{Z}$ , as we now describe.



## 4.2 Efficient Posterior Inference

The posterior updates for  $\mathbf{m}$  and  $\mathbf{u}$  depend on the data only through quantities that can be calculated through  $\mathcal{N}$  and  $\mathcal{P}$ . Let  $n_p(\mathbf{Z}) = \sum_{j: Z_j \leq n_A} I(\gamma_{Z_j, j} = h_p)$  be the number of matching record pairs with agreement pattern  $p$ , and  $N_p = \sum_{j=1}^{n_B} N_{p_j}$  be the total occurrence of pattern  $p$  in the data across all record pairs. Then, conditional on  $\mathbf{Z}$ , we can express the contribution to the likelihood in the full conditional for  $\mathbf{m}$  and  $\mathbf{u}$  as

$$\mathcal{L}(\mathbf{m}, \mathbf{u} \mid \tilde{\gamma}, \mathbf{Z}, \pi) = \prod_{p=1}^P m_p^{n_p(\mathbf{Z})} u_p^{N_p - n_p(\mathbf{Z})}. \quad (15)$$

Additionally, let  $\boldsymbol{\alpha}_0 = (\alpha_{11}, \dots, \alpha_{FL_F})$  and  $\boldsymbol{\beta}_0 = (\beta_{11}, \dots, \beta_{FL_F})$  be concatenated vectors of prior parameters for the  $\mathbf{m}$  and  $\mathbf{u}$  distributions respectively. The terms needed for the posterior updates for the  $\mathbf{m}$  and  $\mathbf{u}$  parameters are given by the appropriate components of the vectors

$$\boldsymbol{\alpha}(\mathbf{Z}^{(s)}) = \boldsymbol{\alpha}_0 + \sum_{p=1}^P n_p(\mathbf{Z}^{(s)}) \times e(h_p), \quad (16a)$$

$$\boldsymbol{\beta}(\mathbf{Z}^{(s)}) = \boldsymbol{\beta}_0 + \sum_{p=1}^P (N_p - n_p(\mathbf{Z}^{(s)})) \times e(h_p). \quad (16b)$$

Specifically, the  $l + \sum_{k=1}^{f-1} L_k$  components of (16a) and (16b) provide the posterior updates for level  $l$  and field  $f$  in (7c) and (7d). However, the vectorized summations can be computed more efficiently than summing over  $n_A n_B$  record pairs for each field and agreement level.

Similarly, the posterior updates for  $\mathbf{Z}$  depend on the data only through quantities calculated through  $\mathcal{N}$ ,  $\mathcal{P}$ , and  $\mathcal{R}$ . For each pattern  $p$ , let  $w_p = m_p/u_p$ . The contribution to the likelihood in the full conditional for  $Z_j$  can be expressed as

$$\mathcal{L}(Z_j \mid \tilde{\gamma}, \mathbf{m}, \mathbf{u}, \pi) = \prod_{p=1}^P u_p^{N_{p_j}} \prod_{i \in r_{p_j}} w_p^{I(Z_j=i)}. \quad (17)$$

With this form, we break this sampling step into two. We first sample among  $P + 1$  options for the agreement pattern between  $B_j$  and its potential link. Define  $r$  as an arbitrary set of records. We have

$$p(Z_j^{(s+1)} \in r \mid \tilde{\gamma}, \mathbf{m}^{(s+1)}, \mathbf{u}^{(s+1)}, \pi^{(s+1)}) \propto \begin{cases} \frac{\pi^{(s+1)} N_{p_j}}{n_A} w_p^{(s+1)}, & r = r_{p_j}; \\ 1 - \pi^{(s+1)}, & r = \{n_A + j\}. \end{cases} \quad (18)$$

Since all records in  $A$  sharing the same agreement pattern with  $B_j$  are equally likely, we then sample among candidate records uniformly using

$$p(Z_j^{(s+1)} = q \mid Z_j^{(s+1)} \in r, \mathbf{m}^{(s+1)}, \mathbf{u}^{(s+1)}, \pi^{(s+1)}) = \begin{cases} \frac{1}{N_{p_j}}, & r = r_{p_j} \text{ and } q \in r; \\ 1, & r = \{n_A + j\} \text{ and } q = n_A + j. \end{cases} \quad (19)$$

## 10 Efficient and Scalable Bipartite Matching with Fast Beta Linkage (*fabl*)

Note that sampling from (21) has complexity  $O(P + 1)$  and sampling (uniformly) from (22) has complexity  $O(1)$ , regardless of the size of  $r_{p_j}$ . In contrast, sampling  $Z_j$  from the full conditional provided in (9) has complexity  $O(n_A)$ , because sampling a value from  $n_A$  options with unequal weights requires normalizing the weights to probabilities.

These changes can greatly improve the speed of the sampler, and each can be parallelized if desired for additional computational speed-ups. We emphasize the computational gains of this split sampler through Lemma 1.

**Lemma 1.** *Let  $n_A$  and  $n_B$  be the number of records in files  $A$  and  $B$ , respectively. Let  $F$  be the number of fields used for comparisons across records, and  $P$  be the number of patterns that comparison vectors exhibit in  $A \times B$ . We assume  $C$  cores available for parallelization and a Gibbs sampler with  $T$  iterations. Then, the overall computational complexity of *fabl* with hashing is  $O(Fn_An_B/C) + O(Tn_BP/C)$ .*

*Proof.* We consider two steps: constructing the comparison vectors and the Gibbs sampler. The computational complexity of all pairwise comparisons across  $A$  and  $B$  is  $O(Fn_An_B)$ . The hashing procedure for all pairwise comparisons is also  $O(Fn_An_B)$ . With  $C$  processors available, we can split these computations across  $C$  equally sized partitions and compute these comparisons in parallel, so the complexity becomes  $O(Fn_An_B/C)$ . There are then trivial computational costs associated with synthesizing summary statistics across these partitions. We note that this contribution to computational complexity applies for all versions of Fellegi and Sunter (1969) algorithms unless they use blocking to reduce the comparison space.

Without hashing, the computational complexity of updating the  $\mathbf{m}$  and  $\mathbf{u}$  parameters is  $O(Fn_An_B)$ . However, by doing calculations over the agreement patterns rather than the individual records, hashing reduces the overall complexity to  $O(P)$ . The complexity of updating  $\mathbf{Z}$  sequentially at the record level as in Sadinle (2017) is  $O(n_An_B)$ . With hashing, we first sample the agreement pattern of the match with complexity  $O(n_BP)$ , and then we sample the record exhibiting that pattern with complexity  $O(n_B)$ . Thus, the complexity of sampling  $\mathbf{Z}$  in a single iteration is  $O(n_BP)$ . Since  $P \ll n_A$  in most applications, we have reduced the complexity of sampling  $\mathbf{Z}$  from  $O(Fn_An_B)$  under BRL to  $O(n_BP)$  under *fabl*. With parallelization, this complexity is further reduced to  $O(n_BP/C)$ , and so the entire Gibbs sampler has complexity  $O(Tn_BP/C)$ . In summary, the total computational complexity for *fabl* is  $O(Fn_An_B/C) + O(Tn_BP/C)$ .  $\square$

### 4.3 Scaling to Large Linkage Tasks

For linkage tasks with large amounts of records, we can partition  $A$  and  $B$  into  $t_A$  and  $t_B$  smaller disjoint batches for more manageable computations. Let  $\{A^1, \dots, A^{t_A}\}$  be a partition of  $A$  such that  $\cup_{a=1}^{t_A} A^a = A$  and  $A^a \cap A^{a'} = \emptyset$  for all  $a \neq a'$ . Likewise, let  $\{B^1, \dots, B^{t_B}\}$  be a partition of  $B$  such that  $\cup_{b=1}^{t_B} B^b = B$  and  $B^b \cap B^{b'} = \emptyset$  for all  $b \neq b'$ . For each  $a$  and  $b$ , we compute comparison vectors for all records in  $A^a \times B^b$  to construct the comparison matrix  $\gamma^{ab}$ .

We then conduct hashing, obtain the compressed  $\tilde{\gamma}^{ab}$ , and delete the memory intensive matrix  $\gamma^{ab}$  before continuing with the next batch of data. In detail, we calculate

$$r_{p_j}^{ab} = \{i \in 1, \dots, n_A \mid \gamma_{ij} = h_p, B_j \in B^b\}, \quad (20a)$$

$$N_{p_j}^{ab} = |r_{p_j}^{ab}|. \quad (20b)$$

These can be computed sequentially or in parallel. Summary statistics from each pairwise batch comparison can be synthesized to recover summary statistics for the full comparison matrix  $\gamma$  through

$$r_{p_j} = \{r_{p_j}^{11}, \dots, r_{p_j}^{t_A t_B}\} \text{ for } a = 1, \dots, t_A \text{ and } b = 1, \dots, t_B, \quad (21a)$$

$$N_{p_j} = \sum_{a=1}^{t_A} \sum_{b=1}^{t_B} N_{p_j}^{ab}. \quad (21b)$$

## 4.4 Storage Efficient Indexing

As discussed in Section 4.1, storing the indices, patterns, and counts in  $\tilde{\gamma}$  uses less memory than storing the full comparison matrix  $\gamma$ . However, recording the indices for all record pairs in  $\mathcal{R}$  can become computationally burdensome for very large linkage tasks. We next introduce storage efficient indexing (SEI), which, when used with the methods in Section 4.3, allows us to compute  $\mathcal{N}$  for all  $n_A n_B$  record pairs while greatly reducing the memory costs of  $\mathcal{R}$  associated with unlikely matches. This allows all-to-all comparisons for substantially larger linkage tasks.

All records  $A_i$  that share agreement pattern  $p$  with record  $B_j$  have the same  $w_p$ . These records have the same probability to be identified as the link for record  $B_j$ . Thus, records  $i \in r_{p_j}$  with large  $N_{p_j}$  are unlikely to be sampled consistently enough to be deemed a match in the Bayes estimate. Rather than store all of these record labels, we store only a small number  $S$ . For each  $r_{p_j}^{ab}$ , we sample  $S$  indices without replacement to form  $\text{SEI}(r_{p_j}^{ab})$ . We collect these memory-reduced lists to form  $\text{SEI}(r_{p_j})$  as in (24a), and collect these to form  $\text{SEI}(\mathcal{R})$ .

As shown in the full conditionals in (16a), (16b), and (21), all original record pairs are still accounted for through  $\mathcal{N}$ . We must revise (22) to account for the reduced number of labels, so instead use

$$p\left(Z_j^{(s+1)} = q \mid Z_j^{(s+1)} \in r, \mathbf{m}^{(s+1)}, \mathbf{u}^{(s+1)}, \pi^{(s+1)}\right) = \begin{cases} \frac{1}{|r|}, & r = \text{SEI}(r_{p_j}) \text{ and} \\ & q \in r; \\ 1, & r = \{n_A + j\} \text{ and} \\ & q = n_A + j \end{cases} \quad (22)$$

in the second step of the split sampler. In this way, we can conduct all posterior inference with the memory reduced  $\text{SEI}(\tilde{\gamma}) = \{\mathcal{P}, \text{SEI}(\mathcal{R}), \mathcal{N}\}$ . We provide guidance on choice of  $S$  through a simulation in Section 5.3.

## 5 Simulation Studies

We demonstrate the speed and accuracy of `fabl` as compared to `BRL` through several simulation studies.

### 5.1 Speed

In our first simulation, we generate comparison vectors from pre-specified distributions so that we can easily increase the size of the linkage problem. We use  $F = 5$  binary comparisons with probabilities for matching and nonmatching pairs shown in Table 1. For each record in  $B$ , we simulate  $n_A$  comparison vectors, resulting in a comparison matrix  $\gamma \in \mathbb{R}^{n_A n_B \times F}$ . For  $n_B/2$  of these records, there is no match in  $A$ , so we simulate  $n_A$  comparison vectors from the  $\mathbf{u}$  probabilities. For the other  $n_B/2$  of these records, there is one match in  $A$ , so we simulate 1 comparison vector from the  $\mathbf{m}$  probabilities, and  $n_A - 1$  comparison vectors from the  $\mathbf{u}$  probabilities. We compare the run-time of `fabl` (with no SEI) against `BRL` as we increase  $n_A$  and  $n_B$ . Since we have five binary comparison fields with no missingness, the number of unique patterns  $P$  is bounded above by  $2^5 = 32$ , a bound which is consistently attained in simulations with more records.

In Figure 1a, where we increase both  $n_A$  and  $n_B$ , `BRL` is faster than `fabl` for low sample sizes, but `fabl` is significantly faster at handling larger sample sizes. In particular, run-time for `BRL` grows quadratically (or linearly with the size of both  $A$  and  $B$ ) while run-time for `fabl` grows linearly (in the size of only  $B$ ).

In Figure 1b, where we fix  $n_B = 500$ , we see near linear growth for the run-time under `BRL` as  $n_A$  increases, and much more static run-time under `fabl`. The slight increases in the run-time for `fabl` are due primarily to the hashing step, which again can be run in parallel for large data. To illustrate that these trends are generalizable to other specifications of the comparison vectors, we have included the run-time results for an additional simulation study, under different comparison vector settings, in Supplement E.

Importantly, `BRL` implements a Gibbs sampler that is coded C (Sadinle, 2017), while `fabl` currently uses non-optimized code written only in R. While this complicates comparisons, and indeed disfavors `fabl`, the computational speed gains for `fabl` are still evident, especially for larger sample sizes. Additionally, although `fabl` is amenable to parallelization, this simulation is run on a single core. Implementing `fabl` in C++ with parallelization for the hashing step and sampling the matching status of the record pairs should lead to even more computational gains.

### 5.2 Accuracy

Computational speed-ups are only worthwhile if not accompanied by a notable loss of record linkage accuracy. Therefore, we examine the accuracy of `fabl` relative to `BRL` by replicating a simulation study from Sadinle (2017). The simulations employ a collection of synthetic data files with varying amounts of error and overlap (the number of records in common across files). Following methods proposed by Christen and Pudjijono (2009)

	$\mathbf{m}$		$\mathbf{u}$	
	Agree	Disagree	Agree	Disagree
First Name	$\frac{19}{20}$	$\frac{1}{20}$	$\frac{1}{100}$	$\frac{99}{100}$
Last Name	$\frac{19}{20}$	$\frac{1}{20}$	$\frac{1}{100}$	$\frac{99}{100}$
Day	$\frac{20}{19}$	$\frac{1}{20}$	$\frac{1}{30}$	$\frac{29}{30}$
Month	$\frac{20}{19}$	$\frac{1}{20}$	$\frac{1}{12}$	$\frac{11}{12}$
Year	$\frac{20}{19}$	$\frac{1}{20}$	$\frac{1}{12}$	$\frac{11}{12}$

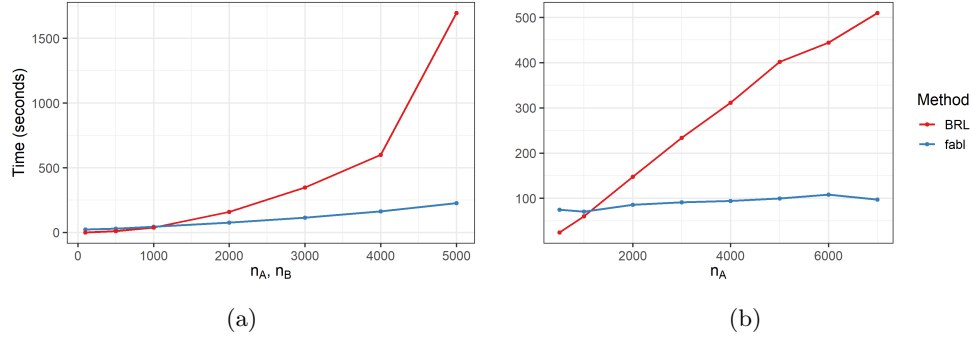
Table 1: Probabilities used for  $\mathbf{m}$  and  $\mathbf{u}$  in simulation study in Section 5.1.

Figure 1: Run-time for BRL and **fabl** to run 1000 Gibbs iterations in simulations described in Section 5.1. In (1a), both  $n_A$  and  $n_B$  are increasing. We see quadratic growth in BRL and linear growth in **fabl**. In (1b), only  $n_A$  only is increasing. We see linear growth in BRL and approximately constant run-time in **fabl**.

and Christen and Vatsalan (2013), clean records are first simulated from frequency tables for first name, last name, age, and occupation in Australia. Fields are then chosen for distortion uniformly at random. Names are subject to string insertions, deletions and substitutions, as well as common keyboard, phonetic, and optical recognition errors. Age and occupation are distorted through keyboard errors and missingness. These synthetic data files are available in the supplement to Sadinle (2017).

We create comparison vectors according to the default settings of the `compareRecords` function from the BRL package, shown in Table 2. Each simulation identifies matched individuals between two data files, each with 500 records. We conduct linkage when matching records exhibit 1, 2, and 3 errors across the four fields, and when there are 50, 250, and 450 individuals in common across data files. Under each of these settings, we use 100 pairs of simulated data files in order to obtain uncertainty quantification on our performance metrics. We use uniform priors for the  $\mathbf{m}$ ,  $\mathbf{u}$ , and  $\pi$  parameters, with  $\alpha_{fl} = \beta_{fl} = 1$  for all  $f$  and  $l$ . We run the Gibbs sampler for 1000 iterations, and discard the first 100 as burn-in. We calculate Bayes estimates  $\hat{\mathbf{Z}}$  of the linkage structure using the loss function and post-processing procedure described in Supplement B. Traceplots for parameters of interest for one example simulation are provided in Supplement C; they show no obvious concern over MCMC convergence. We also replicate this simulation

## 14 *Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)*

Fields	Similarity	Level of Disagreement			
		1	2	3	4
First and Last Name	Levenstein	0	(0, .25]	(.25, .5]	(.5, .1]
Age and Occupation	Binary	Agree	Disagree		

Table 2: Construction of comparison vectors for accuracy study with simulated data files of Section 5.2.

allowing **fabl** to leave some components of the linkage structure undetermined and left for clerical review; those results are in Supplement D.

We compare **fabl** to BRL in terms of recall, precision and F-measure, as defined in Christen (2012). Recall is the proportion of true matches found by the model, that is,  $\sum_{j=1}^{n_B} I(\hat{Z}_j = Z_j, Z_j \leq n_A) / \sum_{j=1}^{n_B} I(Z_j \leq n_A)$ . Precision is the proportion of links found by the model that are true matches, that is,  $\sum_{j=1}^{n_B} I(\hat{Z}_j = Z_j, Z_j \leq n_A) / \sum_{j=1}^{n_B} I(\hat{Z}_j \leq n_A)$ . The F-measure balances the two metrics to provide an overall measure of accuracy, and is defined as  $2(\text{Recall} + \text{Precision}) / (\text{Recall} \times \text{Precision})$ . In Figure 2, we see that the two methods have comparable performance at all levels of error and overlap.

### 5.3 SEI Sensitivity

Finally, our last simulation demonstrates the robustness of **fabl** to different values of  $S$  for the SEI memory reduction procedure. We perform record linkage on one set of the synthetic data files described in Section 5.2 with 500 records in each data file, 250 entities in common across data files, and 3 errors present across matching records. We perform SEI without batching the data, that is,  $t_A = t_B = 1$ . In practice, when it is computationally feasible to create and store  $\gamma$  without batching, there is no need to reduce the memory of the hashed matrix through SEI. For illustration, however, it is easier to examine the effects of choices of  $S$  in this setting.

We perform linkage using SEI with  $S \in (1, 2, 5, 10, 20)$ , and without using SEI, always with 500 iterations of the Gibbs sampler. Any particular SEI implementation may improve or worsen linkage performance; if the SEI procedure happens to only remove pairs that are not matches, recall and precision will improve. Therefore, we perform linkage under each setting 100 times, recording the linkage estimate  $\hat{Z}$ , and recall and precision.

In Figure 3, the largest number of distinct linkage estimates occurs when  $S = 1$ . In this case, the SEI procedure arbitrarily removes large numbers of record labels from consideration, resulting in a noisier estimate of the linkage structure. The number of distinct linkage estimates decreases as  $S$  increases, with larger values of  $S$  providing results more similar to the linkage without SEI. In Figure 4, we see similar patterns in precision. Setting  $S = 1$  can arbitrarily remove the index of a true match, leading the Gibbs sampler to concentrate probability on a false match, while larger values of  $S$  produce results mirroring implementation with no SEI. We note, however, that even with  $S = 1$ , the loss in precision is small in these simulations.

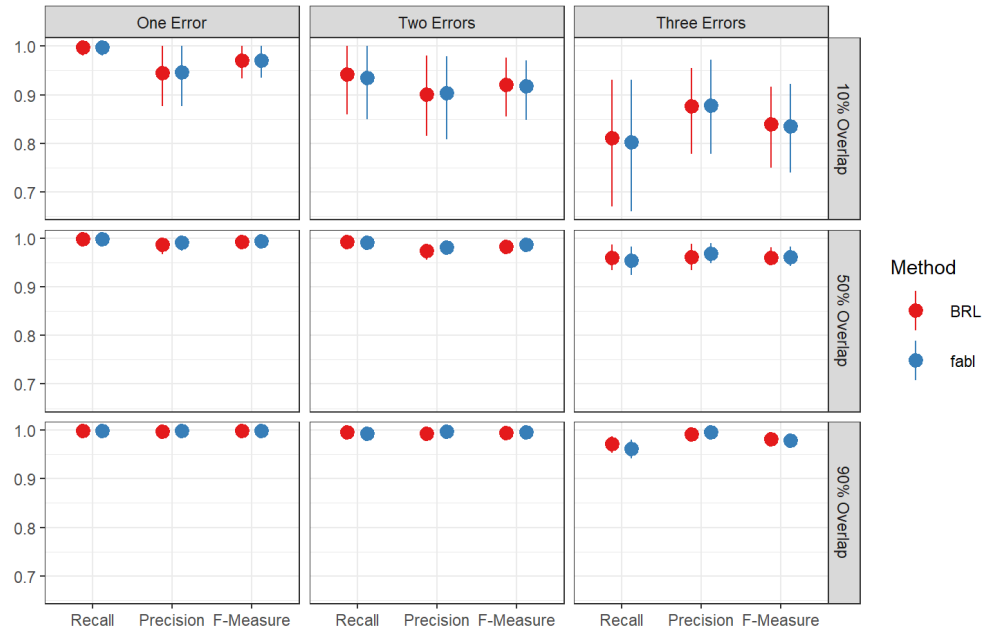


Figure 2: Posterior means and credible intervals for accuracy metrics under the replication of the simulation study from [Sadinle \(2017\)](#). For each level of overlap and each level of error, we have 100 paired sets of 500 records. Thus this table summarizes results for 900 data files. We see comparable performance for all levels of error and overlap.



Figure 3: Distinct values of  $\hat{Z}$  in the simulations of Section 5.3.

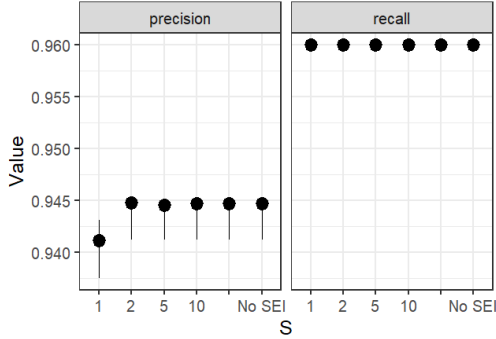


Figure 4: Means and 95% credible intervals for recall and precision in the simulations of Section 5.3.

361 Although the figures suggest that  $S = 2$  is adequate for maintaining linkage per-  
 362 formance, we suggest a more conservative value like  $S = 10$ . When evaluating the  
 363 performance of a record linkage algorithm, researchers often examine posterior proba-  
 364 bilities. By concentrating probability mass on arbitrary nonmatches, low values of  $S$   
 365 may induce artificially high posterior probability for certain record pairs, providing a  
 366 misleading perception of model performance.

## 367 6 Case Studies

368 In our first case study, we revisit data from the El Salvadoran Civil War analyzed by  
 369 [Sadinle \(2017\)](#). Though the data files used in this case study are small, it shows how  
 370 the computational complexity of `fabl` depends on the number of unique agreement  
 371 patterns found in the data, and how significant computational gains can be achieved  
 372 by simplifying the construction of the comparison vectors. In the second case study, we  
 373 apply `fabl` to link records from the National Long Term Care Study, a larger linkage  
 374 task that is not feasible in reasonable time under `BRL` with typical computing setups.

### 375 6.1 Civilian Killings from the El Salvadoran Civil War

376 The country of El Salvador was immersed in civil war from 1980 to 1991. We are  
 377 interested in estimating the total number of individuals killed in the war. We utilize  
 378 lists of documented deaths from the war, one collected by El Rescate - Tutela Regal  
 379 (ERTL) and another from the Salvadoran Human Rights Commission (CDHES, by its  
 380 acronym in Spanish).<sup>1</sup> The ERTL data set comprises digitized denunciations published  
 381 throughout the conflict, and the CDHES data set comprises killings reported directly to  
 382 the organization ([Howland, 2008](#); [Ball, 2000](#); [Green and Ball, 2019](#)). The ERTL required

<sup>1</sup>We thank the Human Rights Data Analysis Group (HRDAG) for granting access to these data.



Fields	Similarity	Level of Disagreement			
		1	2	3	4
First and Last Name	Modified Levenstein	0	(0, .25]	(.25, .5]	(.5, 1]
Year of Death	Absolute Difference	0	1	2	3+
Month of Death	Absolute Difference	0	1	2-3	4+
Day of Death	Absolute Difference	0	1-2	3-7	8+
Municipality	Binary	Agree	Disagree		

Table 3: Construction of comparison vectors for El Salvador data resembling original implementation from [Sadinle \(2017\)](#). This setup leads to 1875 possible agreement patterns in total.

Fields	Similarity	Level of Disagreement		
		1	2	3
First and Last Name	Modified Levenstein	0	(0, .25]	(.25, 1]
Year of Death	Binary	Agree	Disagree	
Month of Death	Binary	Agree	Disagree	
Day of Death	Absolute Difference	0	1	2+
Municipality	Binary	Agree	Disagree	

Table 4: Construction of comparison vectors for El Salvador data for increased speed under **fabl**. This setup leads to 432 possible agreement patterns in total.

383 additional investigation before recording denunciations as human rights abuses, and  
 384 reports to the CDHES were made shortly after the events occurred; thus, both data  
 385 files are thought to be fairly reliable. When estimating the total number of individuals  
 386 killed, one cannot simply sum the numbers recorded by each organization, as it is likely  
 387 that the same individuals are recorded in multiple casualty lists. Instead, record linkage  
 388 techniques must be used to merge data files before analyzing the data ([Lum et al., 2013](#)).

389 There are several challenges with these data. First, the ERTL data file was automat-  
 390 ically digitized, which inherently leads to some degree of typographical error. Second, it  
 391 is common for villages in El Salvador to consist of only four to five extended families.  
 392 This means there are a small number of last names in use, so many individuals have the  
 393 same first and last name. Since the only fields recorded are given name, last name, date  
 394 of death, and place of death, this leads to several instances in which distinct individuals  
 395 have identical records. These individuals may be distant cousins, or are perhaps entirely  
 396 unrelated, but would pose challenges for any record linkage method.

397 Following [Sadinle \(2017\)](#), we utilize records that have non-missing entries for given  
 398 and last name, which results in  $n_A = 4420$  records in CDHES and  $n_B = 1323$  records in  
 399 ERTL. We standardize names to account for common misspellings and use a modified  
 400 Levenstein distance when comparing names to account for the fact that second names  
 401 are often omitted in Spanish. Place of death is recorded by municipality and department  
 402 within that municipality; however, since department is missing in 95% of records in  
 403 CDHES and 80% of records in ERTL, we exclude department from our analysis. Thus,  
 404 we conduct record linkage using given name, last name, municipality, and day, month,  
 405 and year of death. We use uniform priors for the  $\mathbf{m}$ ,  $\mathbf{u}$ , and  $\pi$  parameters.

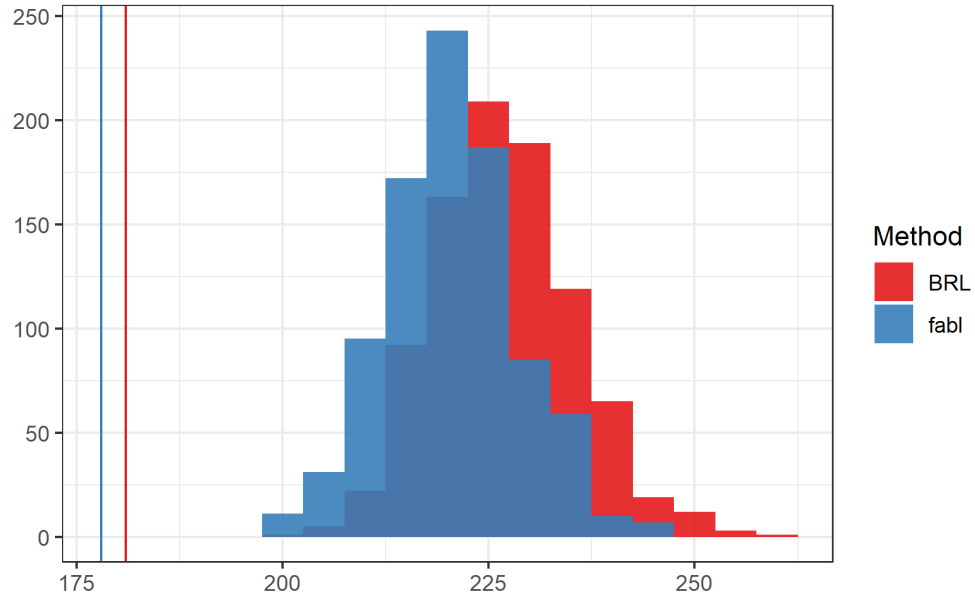


Figure 5: Posterior distribution and Bayes estimate of overlap across the two files in the El Salvador case study. We note they are quite similar under both methods.

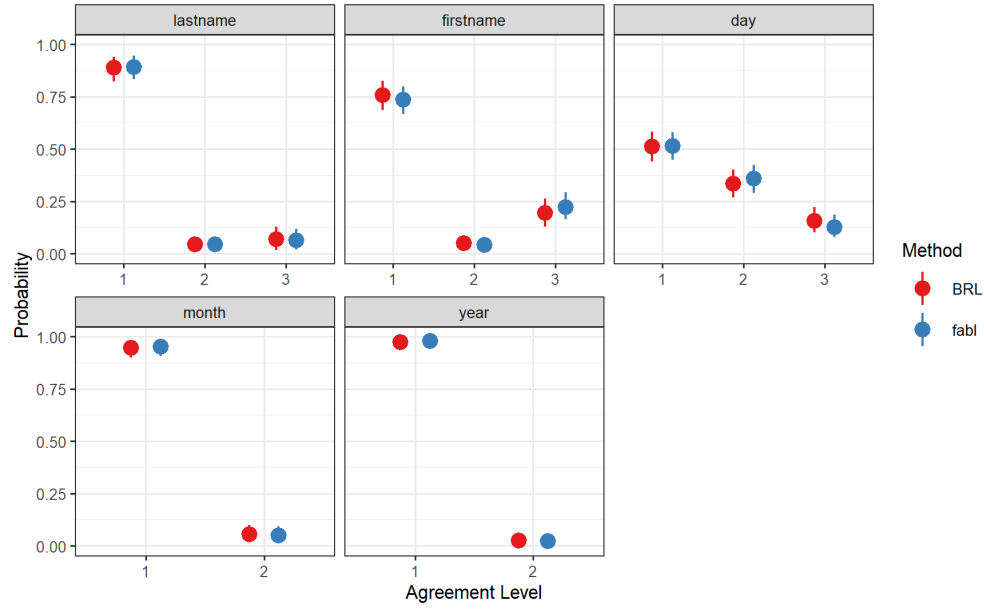


Figure 6: Posterior estimates of  $m$  parameters with 95% credible intervals for the El Salvador case study. They are quite similar across the two methods.

We initially followed the comparison vector constructions set by [Sadinle \(2017\)](#), using four levels of agreement for each field, according to the thresholds provided in [Table 3](#). This results in 1875 possible agreement patterns, with 1173 patterns realized in the data. However, we noticed that the posterior distributions of several levels of the  $\mathbf{m}$  and  $\mathbf{u}$  parameters were nearly identical in an initial run of BRL, suggesting that these levels were unnecessary.

Therefore, we perform our analysis with the agreement levels for each field according to [Table 4](#). Among the 432 possible agreement patterns, 159 are realized in the data. With this revised comparison specification, `fabl` runs in 109 seconds, approximately 3 times faster than the BRL run-time of 313 seconds. The estimates of the  $\mathbf{m}$  parameters under each method are similar, as shown in [Figure 6](#). Estimates of the  $\mathbf{u}$  parameters are indistinguishable, and thus omitted. Traceplots for parameters of interest are provided in [Supplement F](#).

For completeness, we note that linkage with the more detailed comparison vectors requires 945 seconds for BRL, and 1093 seconds for `fabl`. Apparently, the number of patterns is sufficiently many that the computational savings from `fabl` does not overcome the inherent speed differences of C as opposed to R.

Through `fabl`, we arrive at a Bayes estimate of 178 individuals recorded in both data files. We calculate posterior samples of the size of the overlap across files by finding the number of links in each iteration of the Gibbs sampler, and subtracting the number of matches that violate one-to-one matching. The posterior 95% credible interval for the overlap across files is (205, 236), indicating that the Bayes estimate identifies fewer matches than the Gibbs sampler identifies on average. This is because a large number of records in ERTL have multiple plausible matches in CDHES; `fabl` recognizes that a match exists among the several options, but is unable to definitely declare a specific pair as a match in the Bayes estimate. We see similar results under BRL, with a Bayes estimate of 181 individuals recorded in both data files, and a posterior 95% credible interval of (211, 244). See [Figure 5](#) for a visual comparison of the Bayes estimates and posterior credible intervals for the two methods. We note that Bayes estimates falling outside of posterior credible intervals has been observed previously in the record linkage literature ([Sadinle, 2017](#); [Steorts et al., 2016](#)), and remains a topic for future research.

## 6.2 National Long Term Care Study

The National Long Term Care Study (NLTCS) is a longitudinal study tracking the health outcomes of Medicare recipients ([Steorts et al., 2016](#)). The initial survey began in 1982, with follow-up surveys taken approximately every five years. As such, patients are surveyed at most once in a given year, and many patients are surveyed across multiple years. In addition, patients can either drop out of the study, pass away, or enter as new patients. Hence, the assumptions of our model hold for this study. We seek to link records over the  $n_A = 20485$  individuals from 1982 to the  $n_B = 17466$  individuals from 1989. The NLTCS data have longitudinal links, so that in reality one does not need to conduct record linkage. However, following the strategy in [Guha et al. \(2022\)](#), we break the longitudinal links and treat the data from 1982 and 1989 as stand-alone data files.

Fields	Similarity	Level of Disagreement		
		1	2	3
Sex	Binary	Agree	Disagree	
Year of Birth	Binary	Agree	Disagree	
Month of Birth	Binary	Agree	Disagree	
Day of Birth	Binary	Agree	Disagree	
Location	Custom	Same State and Office	Same State	Otherwise

Table 5: Construction of comparison vectors for NLTCs data.

We link records using sex, location, and day, month, and year of birth using the criteria shown in Table 5. Storing  $\gamma$  constructed through three comparison scores for each of the approximately 400 million record pairs would require around 8 GB of memory. Standard settings on a 16 GB personal computer do not allow storage of an object of this size, and thus BRL is unable to perform this linkage task on such a machine. However, through the method described in Section 4.3, we perform 30 smaller comparison tasks, using  $t_A = 1$  and  $t_B = 30$ . We conduct linkage with all record indices recorded and also with SEI using  $S = 10$ , and obtain identical results. The hashed  $\tilde{\gamma}$  without SEI is about 2.2 GB, and with SEI, it is about 760 MB. Constructing the comparisons sequentially took approximately 40 minutes, which could be reduced considerably through parallel computing.

We run a Gibbs sampler for 1000 iterations, taking about 235 seconds. Traceplots do not suggest convergence issues, and are similar to those seen in Supplement C and F. As shown in Figure 7, the Bayes estimate of the linkage structure has 9634 matches, with a 95% credible interval of (9581, 9740). Since we have access to the true linkage structure, we can calculate recall to be 0.89 and precision to be 0.98, resulting in an F-measure of 0.94.

## 7 Conclusion

In this paper, we have proposed **fabl**, a Bayesian record linkage method that extends the work of Sadinle (2017) to scale to large data sets. We have proven that the proposed hashing method and model assumptions allow for a linkage procedure whose computational complexity does not scale with the size of the larger data file. This makes **fabl** computationally advantageous in many linkage scenarios, particularly when one data file is substantially smaller than the other. We have also shown that storage efficient indexing, in tandem with hashing, greatly reduces the memory costs required for all-to-all comparisons, giving practitioners an option for larger record linkage tasks potentially even without the use of blocking or indexing. We have demonstrated the speed and accuracy of **fabl** by replicating a simulation study and a case study in Sadinle (2017), and through an additional case study that is computationally impractical under BRL.

Although the **fabl** method greatly reduces the memory costs for all-to-all comparisons, computing the comparisons for all  $n_A n_B$  record pairs still can be prohibitive for larger linkage tasks. Indeed, constructing the comparison vectors for the NLTCs linkage task

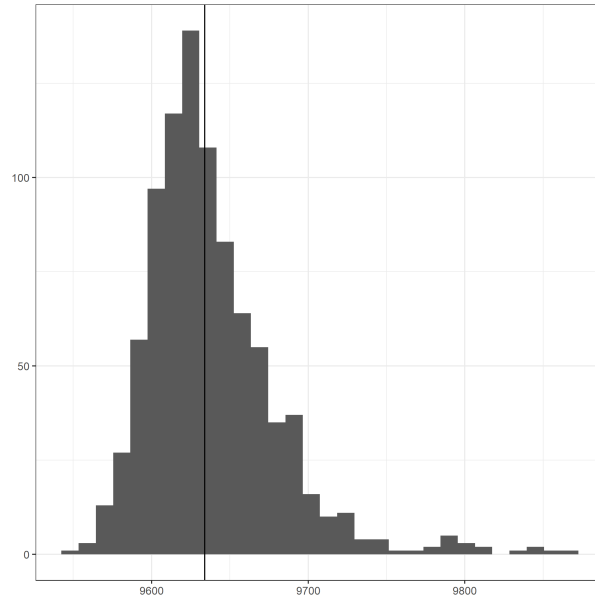


Figure 7: Posterior distribution and Bayes estimate of overlap across years 1982 and 1989 of NLTCs data.

involving around 40,000 records in Section 6.2 took around 40 minutes. Due to the quadratic nature of the comparison space, this computation time would grow quickly with the size of the linkage task, and would be extremely slow when dealing with millions of records. Although it is common to use deterministic blocking to reduce the comparison space and then apply probabilistic record linkage within each block, issues arise when sizes of blocks vary across the linkage task. In future work, we seek to extend `fabl` to account for such deterministic blocking, making the framework amenable to even larger linkage tasks.

**Acknowledgements:** We thank the reviewers for extensive comments and suggestions that immensely improved the quality of the article.

## Supplementary Material

Supplementary Material of “Efficient and Scalable Bipartite Matching with Fast Beta Linkage (`fabl`).” See the .pdf document.

## References

- Aleshin-Guendel, S. and Sadinle, M. (2023), “Multifile Partitioning for Record Linkage and Duplicate Detection,” *Journal of the American Statistical Association*, 0, 1–10. [1](#)
- Ball, P. (2000), “The Salvadoran Human Rights Commission: Data Processing, Data Representation, and Generating Analytical Reports,” in *Making the Case: Investigating Large Scale Human Rights Violations Using Information Systems and Data Analysis*, eds. P. Ball, H. F. Spierer, and L. Spierer, pp. 15–24, American Association for the Advancement of Science. [17](#)
- Betancourt, B., Sosa, J., and Rodríguez, A. (2022), “A prior for record linkage based on allelic partitions,” *Computational Statistics & Data Analysis*, 172, 107 – 474. [1](#)
- Bilenko, M. and Mooney, R. (2006), “Riddle: Repository of Information on Duplicate Detection, Record Linkage, and Identity Uncertainty,” Online; retrieved July 29, 2020. [3](#)
- Christen, P. (2012), “A Survey of Indexing Techniques for Scalable Record Linkage and Deduplication,” *IEEE Transactions on Knowledge and Data Engineering*, 24, 1537–1555. [1](#), [14](#)
- Christen, P. (2019), “Data Linkage: The Big Picture,” *Harvard Data Science Review*, 1, <https://hdsr.mitpress.mit.edu/pub/8fm8lo1e>. [7](#)
- Christen, P. and Pudjijono, A. (2009), “Accurate Synthetic Generation of Realistic Personal Information,” in *Advances in Knowledge Discovery and Data Mining*, eds. T. Theeramunkong, B. Kijirikul, N. Cercone, and T.-B. Ho, pp. 507–514, Berlin, Heidelberg, Springer Berlin Heidelberg. [14](#)
- Christen, P. and Vatsalan, D. (2013), “Flexible and Extensible Generation and Corruption of Personal Data,” in *Proceedings of the 22nd ACM International Conference on Information and Knowledge Management, CIKM ’13*, p. 1165–1168, New York, NY, USA, Association for Computing Machinery. [14](#)
- Cohen, W. W., Ravikumar, P., and Fienberg, S. E. (2003), “A Comparison of String Distance Metrics for Name-Matching Tasks,” in *Proceedings of the 2003 International Conference on Information Integration on the Web*, p. 73–78, AAAI Press. [3](#)
- Dalzell, N. M. and Reiter, J. P. (2018), “Regression Modeling and File Matching Using Possibly Erroneous Matching Variables,” *Journal of Computational and Graphical Statistics*, 27, 728–738. [1](#)
- Elmagarmid, A. K., Ipeirotis, P. G., and Verykios, V. S. (2007), “Duplicate Record Detection: A Survey,” *IEEE Transactions on Knowledge and Data Engineering*, 19, 1–16. [3](#)
- Enamorado, T., Fifield, B., and Imai, K. (2019), “Using a Probabilistic Model to Assist Merging of Large-Scale Administrative Records,” *American Political Science Review*, 113, 353–371. [1](#), [4](#), [7](#)
- Fair, M. (2004), “Generalized Record Linkage System—Statistics Canada’s Record Linkage Software,” *Austrian Journal of Statistics*, 33, 37–53. [1](#)

- 533 Fellegi, I. P. and Sunter, A. B. (1969), “A Theory for Record Linkage,” *Journal of the*  
 534 *American Statistical Association*, 64, 1183–1210. [1](#), [2](#), [3](#), [4](#), [5](#), [7](#), [11](#)
- 535 Fortunato, S. (2010), “Community Detection in Graphs,” *Physics Reports*, 486, 75–174.  
 536 [2](#)
- 537 Gill, L. and Goldacre, M. (2003), “English National Record Linkage of Hospital Episode  
 538 Statistics and Death Registration Records,” *Report to the Department of Health*. [1](#)
- 539 Green, A. H. and Ball, P. (2019), “Civilian Killings and Disappearances During Civil  
 540 War in El Salvador (1980–1992),” *Demographic Research*, 41, 781–814. [17](#)
- 541 Guha, S., Reiter, J., and Mercatanti, A. (2022), “Bayesian Causal Inference with Bipartite  
 542 Record Linkage,” *Bayesian Analysis*, 17, 1275–1299. [20](#)
- 543 Gutman, R., Afendulis, C. C., and Zaslavsky, A. M. (2013), “A Bayesian Procedure  
 544 for File Linking to Analyze End-of-Life Medical Costs,” *Journal of the American*  
 545 *Statistical Association*, 108, 34–47. [1](#)
- 546 Howland, T. (2008), “How El Rescate, a Small Nongovernmental Organization, Con-  
 547 tributed to the Transformation of the Human Rights Situation in El Salvador,” *Human*  
 548 *Rights Quarterly*, 30, 703–757. [17](#)
- 549 Jaro, M. A. (1989), “Advances in Record-Linkage Methodology as Applied to Matching  
 550 the 1985 Census of Tampa, Florida,” *Journal of the American Statistical Association*,  
 551 84, 414–420. [2](#), [4](#)
- 552 Larsen, M. D. (2005), “Advances in Record Linkage Theory: Hierarchical Bayesian  
 553 Record Linkage Theory,” in *Proceedings of the Joint Statistical Meetings, Section on*  
 554 *Survey Research Methods*, pp. 3277–3284, The American Statistical Association. [2](#)
- 555 Larsen, M. D. and Rubin, D. B. (2001), “Iterative Automated Record Linkage Using  
 556 Mixture Models,” *Journal of the American Statistical Association*, 96, 32–41. [4](#)
- 557 Little, R. and Rubin, D. (2002), *Statistical Analysis with Missing Data*, Wiley, Hoboken,  
 558 New Jersey. [3](#)
- 559 Lum, K., Price, M. E., and Banks, D. (2013), “Applications of Multiple Systems Estima-  
 560 tion in Human Rights Research,” *The American Statistician*, 67, 191–200. [17](#)
- 561 Marchant, N. G., Kaplan, A., Elazar, D. N., Rubinstein, B. I. P., and Steorts, R. C.  
 562 (2021), “d-blink: Distributed End-to-End Bayesian Entity Resolution,” *Journal of*  
 563 *Computational and Graphical Statistics*, 30, 406–421. [1](#)
- 564 Murray, J. S. (2016), “Probabilistic Record Linkage and Deduplication after Indexing,  
 565 Blocking, and Filtering,” *Journal of Privacy and Confidentiality*, 7, 3–24. [7](#)
- 566 Newcombe, H. B., Kennedy, J. M., Axford, S. J., and James, A. P. (1959), “Automatic  
 567 Linkage of Vital Records,” *Science*, 130, 954–959. [1](#)
- 568 Sadinle, M. (2017), “Bayesian Estimation of Bipartite Matchings for Record Linkage,”  
 569 *Journal of the American Statistical Association*, 112, 600–612. [1](#), [2](#), [4](#), [5](#), [6](#), [11](#), [13](#), [14](#),  
 570 [15](#), [17](#), [19](#), [20](#), [21](#)

- 571 Steorts, R. C., Ventura, S. L., Sadinle, M., and Fienberg, S. E. (2014), “A Comparison  
572 of Blocking Methods for Record Linkage,” in *Privacy in Statistical Databases*, ed.  
573 J. Domingo-Ferrer, pp. 253–268, Cham, Springer International Publishing. [7](#)
- 574 Steorts, R. C., Hall, R., and Fienberg, S. E. (2016), “A Bayesian Approach to Graphical  
575 Record Linkage and Deduplication,” *Journal of the American Statistical Association*,  
576 111, 1660–1672. [1](#), [20](#)
- 577 Tancredi, A., Liseo, B., et al. (2011), “A Hierarchical Bayesian Approach to Record  
578 Linkage and Population Size Problems,” *The Annals of Applied Statistics*, 5, 1553–1585.  
579 [1](#)
- 580 Tang, J., Reiter, J. P., and Steorts, R. C. (2020), “Bayesian Modeling for Simultaneous  
581 Regression and Record Linkage,” in *Privacy in Statistical Databases*, eds. J. Domingo-  
582 Ferrer and K. Muralidhar, pp. 209–223, Cham, Springer International Publishing.  
583 [1](#)
- 584 Wagner, D., Lane, M., et al. (2014), “The Person Identification Validation System (PVS):  
585 Applying the Center for Administrative Records Research and Applications’ (CARRA)  
586 Record Linkage Software,” Tech. rep., Center for Economic Studies, U.S. Census  
587 Bureau. [1](#)
- 588 Winkler, W. and Thibaudeau, Y. (1990), “An Application of the Fellegi-Sunter Model  
589 of Record Linkage to the 1990 US Decennial Census,” *U.S. Census Research Report*,  
590 pp. 1–22. [1](#)
- 591 Winkler, W. E. (1999), “The State of Record Linkage and Current Research Problems,”  
592 Tech. rep., Statistical Research Division, U.S. Bureau of the Census. [4](#)
- 593 Wortman, J. P. H. (2019), “Record linkage methods with applications to causal inference  
594 and election voting data,” Ph.D. thesis, Duke University. [5](#)