

Original citation:

Ahlstrom, Christina, Muellner, Petra, Spencer, Simon E. F., Hong, Samuel, Saupe, Amy, Rovira, Albert, Hedberg, Craig, Perez, Andres, Muellner, Ulrich and Alvarez, Julio. (2017) Inferring source attribution from a multi-year multi-source dataset of Salmonella in Minnesota. *Zoonoses and Public Health*.

Permanent WRAP URL:

<http://wrap.warwick.ac.uk/86165>

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions. Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

"This is the peer reviewed version of Ahlstrom, Christina, Muellner, Petra, Spencer, Simon E. F., Hong, Samuel, Saupe, Amy, Rovira, Albert, Hedberg, Craig, Perez, Andres, Muellner, Ulrich and Alvarez, Julio. (2017) Inferring source attribution from a multi-year multi-source dataset of Salmonella in Minnesota. *Zoonoses and Public Health*. which has been published in final form at <http://doi.org/10.1111/zph.12351> . This article may be used for non-commercial purposes in accordance with [Wiley Terms and Conditions for Self-Archiving](#)."

A note on versions:

The version presented here may differ from the published version or, version of record, if you wish to cite this item you are advised to consult the publisher's version. Please see the 'permanent WRAP URL' above for details on accessing the published version and note that access may require a subscription.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk



Inferring source attribution from a multi-year multi-source dataset of Salmonella in Minnesota

Journal:	<i>Zoonoses and Public Health</i>
Manuscript ID	ZPH-Oct-16-318.R1
Manuscript Type:	Original Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Ahlstrom, Christina; Epi-interactive, Muellner, P; Epi-interactive Spencer, Simon; University of Warwick Hong, Samuel; University of Minnesota, Veterinary Population Medicine Saupe, Amy; Minnesota Department of Health Rovira, Albert; University of Minnesota, Veterinary Diagnostic Laboratory Hedberg, Craig; University of Minnesota, Division of Environmental Health Sciences Perez, Andres; University of Minnesota, Veterinary Population Medicine Muellner, Ulrich; Epi-interactive Alvarez, Julio; University of Minnesota, Veterinary Population Medicine
Subject Area:	Salmonella spp, Molecular epidemiology, Source attribution, Data visualization

SCHOLARONE™
Manuscripts

1
2
3 1 Inferring source attribution from a multi-year multi-source dataset of *Salmonella* in
4
5 2 Minnesota
6
7
8 3
9
10
11 4 Christina Ahlstrom¹, Petra Muellner¹, Simon EF Spencer², Samuel Hong³, Amy Saupe⁴,
12
13 5 Albert Rovira⁵, Craig Hedberg⁶, Andres Perez³, Ulrich Muellner¹, Julio Alvarez³
14
15
16
17 6
18
19
20 7 ¹ Epi-interactive, Wellington, New Zealand
21
22
23 8 ² University of Warwick, Coventry, United Kingdom
24
25
26 9 ³ Department of Veterinary Population Medicine, College of Veterinary Medicine, University
27
28 10 of Minnesota, St Paul, USA
29
30
31 11 ⁴ Minnesota Department of Health, Saint Paul, USA
32
33
34
35 12 ⁵ Veterinary Diagnostic Laboratory, College of Veterinary Medicine, University of
36
37 13 Minnesota, St Paul, USA
38
39
40 14 ⁶ Division of Environmental Health Sciences, School of Public Health, University of
41
42 15 Minnesota, Minneapolis, USA
43
44
45
46 16 Corresponding author: Petra Muellner, Epi-interactive, Wellington, New Zealand.
47
48 17 petra@epi-interactive.com
49
50
51 18
52
53
54
55
56
57
58
59
60

19 Summary

20 *Salmonella enterica* is a global health concern because of its widespread association with
21 foodborne illness. Bayesian models have been developed to attribute the burden of human
22 salmonellosis to specific sources with the ultimate objective of prioritizing intervention
23 strategies. Important considerations of source attribution models include the evaluation of the
24 quality of input data, assessment of whether attribution results logically reflect the data
25 trends, and identification of patterns within the data that might explain the detailed
26 contribution of different sources to the disease burden. Here, more than 12,000 non-typhoidal
27 *Salmonella* isolates from human, bovine, porcine, chicken, and turkey sources that originated
28 in Minnesota were analyzed. A modified Bayesian source attribution model (available in a
29 dedicated R package), accounting for non-sampled sources of infection, attributed 4,672
30 human cases to sources assessed here. Most (60%) cases were attributed to chicken, though
31 there was a spike in cases attributed to a non-sampled source in the second half of the study
32 period. Molecular epidemiological analysis methods were used to supplement risk modelling
33 and a visual attribution application was developed to facilitate data exploration and
34 comprehension of the large multi-year dataset assessed here. A large amount of within-source
35 diversity and low similarity between sources was observed and visual exploration of data
36 provided clues into variations driving the attribution modelling results. Results from this
37 pillared approach provided first attribution estimates for *Salmonella* in Minnesota and offer
38 an understanding of current data gaps as well as key pathogen population features, such as
39 serotype frequency, similarity and diversity across the sources. Results here will be used to
40 inform policy and management strategies ultimately intended to prevent and control
41 *Salmonella* infection in the state.

1
2
3 42 **Keywords:** Source attribution, molecular epidemiology, data visualization, *Salmonella*,
4
5 43 salmonellosis
6

7 44
8

9
10 45 **Bullet points:**

- 11
12 46 • We demonstrate the attribution of salmonellosis to different sources can be enhanced
13
14 47 through a pillared approach that incorporates data visualization and molecular
15
16 48 epidemiology on top of source attribution modelling. In this study, such
17
18 49 supplementary analyses supported findings from the attribution model, demonstrating
19
20 50 for example high genotype diversity and low similarity between sources.
21
22
23 51 • We developed a modified Bayesian source attribution model that attributed the
24
25 52 majority of salmonellosis cases in Minnesota to chickens. Accounting for a known
26
27 53 data gap, we were able to demonstrate the potential impact of a non-sampled source
28
29 54 on the number of attributed cases.
30
31
32 55 • A visual attribution application enabled users to dynamically explore the occurrence
33
34 56 of *Salmonella* serotypes in different sources over time. The ability to interact and
35
36 57 filter the data assisted in the detection of data irregularities, facilitating accurate
37
38 58 attribution estimates and the interpretation of results.
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

59 Introduction

60 Non-typhoidal *Salmonella enterica* is considered one of the leading causes of foodborne
61 illness in the United States (Scallan et al., 2011) and other countries (Kirk et al., 2015),
62 despite large-scale control efforts initiated by industry, government and consumers (Crim et
63 al., 2015). In the United States, *S. enterica* has the highest reported incidence in humans
64 among bacterial pathogens, with 15.3 cases per 100,000 persons in 2014 (CDC, 2016), the
65 majority of which are sporadic with unknown source of origin (Batz et al., 2005; Ebel et al.,
66 2016). In Minnesota, the incidence of culture-confirmed *Salmonella* cases in 2015 was 17.9
67 per 100,000 persons (Minnesota Department of Health, 2015). The most common form of
68 clinical illness associated with *S. enterica* infection is gastroenteritis, which is mostly non-
69 severe and self-limiting; however, *S. enterica* can cause severe disease or complications in
70 some patients (e.g. sepsis). *S. enterica* has multiple reservoirs, including livestock and
71 domestic pets (Kingsley and Bäumlner, 2000), impairing efforts to identify infection sources
72 and transmission routes. Whereas some *S. enterica* serotypes have adapted to individual host
73 species, others are frequently isolated from a broad range of species and environments
74 (Uzzau et al., 2000).

75 Source attribution models aim to estimate the proportion of human cases attributable to
76 specific sources. Estimates obtained can subsequently be used to guide risk managers and
77 decision-makers in the design and implementation of effective intervention strategies (Pires
78 and Hald, 2010; Sears et al., 2011). In the United States, different source attribution methods
79 using a variety of data sources have been employed to estimate the burden of human
80 salmonellosis from several food commodities. For example, the relative proportions of
81 domestically acquired, sporadic human *Salmonella* infections between 1998 and 2003 among
82 multiple food sources was investigated using a molecular subtyping approach (Guo et al.,

1
2
3 83 2011), whereas other source attribution assessments relied either on expert elicitation
4
5 84 (Hoffmann et al., 2007) or outbreak-associated illnesses (Batz et al., 2012; Gould et al., 2013;
6
7 85 Painter et al., 2013). Data collected during outbreak investigations are typically more
8
9
10 86 accessible than sporadic illness data; however, foodborne outbreaks account only for a small
11
12 87 proportion of total *Salmonella* infections in humans and may not be representative of the
13
14 88 majority of human cases of salmonellosis (Ebel et al., 2016; Painter et al., 2013).

15
16
17 89 Although early risk assessment models estimated the contribution of a single source to the
18
19 90 burden of human illness, Bayesian source attribution models for salmonellosis (Hald et al.,
20
21 91 2004) were subsequently developed to consider multiple sources simultaneously while
22
23 92 allowing for uncertainty around input parameters. Such an approach has been applied
24
25 93 internationally and modified over time to accommodate different contexts and pathogens
26
27 94 (Mullner et al. 2009a; Little et al. 2010; Mughini-Gras et al. 2014; David et al. 2013; Guo et
28
29 95 al. 2011; Glass et al. 2015; Mullner et al. 2009b). Briefly, the model framework compares the
30
31 96 distribution of pathogen serotypes in the source populations to the distribution of serotypes
32
33 97 observed in humans via a Poisson regression model fitted within a Bayesian framework.
34
35
36 98 While models developed thus far have provided valuable insights into the contribution of
37
38 99 different sources to the human disease burden, they generally rely on a specific type of data
39
40 100 that may not always be available. Previous work has, for example, explored the use of
41
42 101 different types of prevalence data (Mullner et al. 2009a). The work presented here proposes a
43
44 102 further extension that may allow the model to be successfully used in the presence of data
45
46 103 gaps.

47
48
49
50
51 104 Bayesian source attribution models often rely on multi-level data collected in multiple years
52
53 105 and sources, requiring significant resources to manually clean and prepare the data. Model
54
55 106 outputs are typically one-dimensional and primarily focus on one quantitative figure (i.e. the
56
57
58
59
60

1
2
3 107 attributed proportion of human cases to each source). A comprehensive understanding of the
4
5 108 characteristics of the input data and results could enhance the way attribution results are
6
7 109 understood by, and communicated to, risk managers and other stakeholders. For example,
8
9 110 molecular epidemiological analyses, such as diversity statistics and similarity indexes
10
11 111 (Muellner et al., 2011), can be incorporated into the analysis to support model outputs.
12
13 112 Additionally, the integration of visual data exploration creates opportunities to dynamically
14
15 113 interact with layered data and identify trends to help explain attribution results. Such an
16
17 114 approach can help make large, complex datasets digestible through interactive graphics that
18
19 115 facilitate incremental exploration of the data (Carroll et al., 2014).
20
21
22
23

24 116 The objective of the study presented here was to explore the feasibility of using a pillared
25
26 117 attribution approach, adding molecular epidemiology and visual data analysis, to a
27
28 118 customised Hald Model based on the work by Mullner et al. (2009), to generate preliminary
29
30 119 estimates of the source-specific burden of salmonellosis in Minnesota over a ten-year period
31
32 120 using state-level data. Results here will help to inform policy and management activities
33
34 121 intended to prevent and control the disease in the state. Additionally, methods presented here
35
36 122 may serve as a framework that could be applied to the attribution of sources of infection for
37
38 123 *Salmonella* and other foodborne pathogens in the United States and other regions.
39
40
41
42

43 124

46 125 Materials and methods

49 126 Data sources

50
51 127 Information from human salmonellosis cases reported to the Minnesota Department of Health
52
53 128 (MDH) from 2005 to 2014 were collected. Data included case serotype (determined by the
54
55 129 MDH Public Health Laboratory [PHL]), date of specimen collection, international travel in
56
57
58
59
60

1
2
3 130 the seven days prior to illness onset (self-reported by cases during routine exposure
4
5 131 interviews at MDH), and whether the case was part of an identified outbreak. Consistent with
6
7 132 the approach developed by Hald et al (2004), cases with a history of international travel were
8
9 133 excluded, as were cases attributed to an outbreak.
10
11
12 134 Information on *Salmonella* isolates of food animal (cattle, swine, poultry) origin isolated by
13
14 135 the Minnesota Veterinary Diagnostic Laboratory (MVDL) from diagnostic submissions
15
16 136 between 2006 and 2015 were also collected. A detailed description of the MVDL database
17
18 137 and cattle and swine isolates is available elsewhere (Hong et al., 2016). Finally, data on the
19
20 138 relative frequency of *Salmonella* serotypes stratified per meat product (chicken, ground
21
22 139 turkey, ground beef, pork chop) from various retail locations in Minnesota that were collected
23
24 140 as part of the FoodNet/National Antimicrobial Resistance Monitoring System (NARMS)
25
26 141 Retail Food Study between 2002 and 2013 were provided by MDH. These food commodity
27
28 142 source data represented 4% of isolates derived from non-human sources and were combined
29
30 143 with the animal data (e.g. isolates from pork chops were added to animal porcine isolates).
31
32 144 The majority (96%) of food-derived isolates came from chicken breasts and ground turkey.
33
34 145 Non-human source categories were limited to bovine, porcine, chicken, and turkey, because
35
36 146 typed isolates from other sources were considered too scarce ($n < 50$) to be included in the
37
38 147 risk model. Further, no data on eggs were available.
39
40
41
42
43
44 148 The consistency of serotype naming was checked within and between data sources. All
45
46 149 serotypes were defined using naming conventions proposed by the Pasteur Institute (Grimont
47
48 150 and Weill, 2007). After 2012, the Centers for Disease Control and Prevention (CDC)
49
50 151 recommended naming all *S. Typhimurium* var. 5 – (formerly var. Copenhagen) as *S.*
51
52 152 *Typhimurium* (*S.* I 4,[5],12:i:1,2; CDC, 2014). Thus, the MDH did not electronically record
53
54 153 variants of *S. Typhimurium* after 2011; however, variants continued to be identified and
55
56
57
58
59
60

1
2
3 154 recorded in paper records. Therefore, *S. Typhimurium* var. 5– (formerly var. Copenhagen)
4
5 155 cases were retrospectively recorded electronically in the current dataset.
6
7

8 156 Diversity and similarity statistics 9

10
11 157 The dataset containing all *Salmonella* isolates from bovine, porcine, chicken, turkey, and
12
13 158 sporadic and domestic human cases was used to calculate diversity and similarity statistics.
14
15 159 Serotype richness and the diversity of *Salmonella* serotypes from different sources was
16
17 160 estimated using rarefaction and Simpson’s index of diversity. The function RAREFY in the
18
19 161 package VEGAN (version 2.2-3) was implemented in R (version 3.2.3; R Core Team 2016)
20
21 162 and Simpson’s index of diversity was calculated in Past (version 3.10; Hammer et al. 2009),
22
23 163 with 9,999 bootstrap replicates. The similarity of serotypes between different sources was
24
25 164 investigated by calculating the proportional similarity index (PSI), which measures the area
26
27 165 of overlap between two frequency distributions of each serotype between sources. Bootstrap
28
29 166 confidence intervals were estimated as described by Mullner et al. (2009).
30
31
32
33

34 167 Source attribution modelling 35

36 168 Data preparation 37 38 39

40 169 *Salmonella* isolates originating from non-human sources that did not match a serotype found
41
42 170 in Minnesota human cases during the study period were excluded. Similarly, *Salmonella*
43
44 171 isolates originating from humans that did not match a serotype found in non-human sources
45
46 172 were excluded. Serotypes that included fewer than five human cases were combined into a
47
48 173 serotype designated as “other”, similar to the approach taken in previous studies (Hald et al.,
49
50 174 2004; Mullner et al., 2009a; David et al., 2012). Sampling effort of non-human sources varied
51
52 175 across the study period; thus, isolates from non-human sources were not segregated into
53
54 176 years. Human isolates were summarized as total cases per serotype per year between 2005-
55
56 177 2014.
57
58
59
60

1
2
3 178 Attribution model

4
5 179 The source attribution model described by Hald et al. (2004) and modified by Muellner et al.
6
7
8 180 (2009) was further modified here to simplify and improve applicability and interpretability of
9
10 181 results. The approach proposed here may be used in a wider variety of circumstances
11
12 182 compared to earlier methods, in particular, without knowledge of the absolute prevalence in
13
14 183 each source.

15
16
17
18 184 The Poisson regression structure was retained (adapted slightly to include time)

19
20 185 $Y_{it} \sim Pois(\sum_{j=1}^S \lambda_{ij})$, where Y_{it} is the number of serotype i human cases in year $t = 1, \dots, T$.

21
22 186 The equation defining the expected number of cases of serotype i (for $i=1, \dots, I$) from source j
23
24 187 (for $j=1, \dots, S$) was modified as

25
26
27
28 188
$$\lambda_{ij} = r_{ij} q_i a_j.$$

29
30
31 189 Model parameters are described in Table 1. Note that this equation does not require the
32
33 190 absolute source prevalences p_{ij} or π_j (Mullner et al., 2009a) and instead uses just the relative
34
35 191 prevalences r_{ij} , plus the strain i - and the source j -specific factors q_i and a_j , respectively. This
36
37 192 model can therefore be fitted when source prevalence data are not available. Let X_{ij} denote
38
39 193 the number of serotype i isolates observed in source j . The following prior distributions were
40
41 194 assumed:

42
43
44
45
46 195 $(r_{1j}, r_{2j}, \dots, r_{Ij}) \sim \text{Dirichlet}(\gamma_1 + X_{1j}, \gamma_2 + X_{2j}, \dots, \gamma_I + X_{Ij})$

47
48
49
50 196 $q_i \sim \text{Gamma}(\theta, \theta)$, with $\theta \sim \text{Gamma}(\alpha^{(\theta)}, \beta^{(\theta)})$

51
52
53
54 197 $a_j \sim \text{Exp}(\alpha_j)$

55
56
57
58 198 The hyperparameter θ represents the precision and also the shape parameter of the

1
2
3 199 distribution of the random effects describing the strain-specific differences q_i . In previous
4
5 200 versions of the source attribution model, the prior for the source-specific parameter a_j was
6
7 201 difficult to specify. Here, because the prior mean of q_i is one and the relative prevalences sum
8
9
10 202 up to one, the prior mean of a_j (which equals α_j^{-1}) is the expected number of cases from
11
12 203 source j . This interpretation allows informed priors for the source-specific factors to be easily
13
14 204 defined.

15
16
17
18 205 In the current implementation, we chose $\alpha_j^{-1} = \sum_{i=1}^I \sum_{t=1}^T Y_{it} / TS$ for all $j = 1, \dots, S$, which
19
20 206 implies the prior belief that an equal number of cases comes from each source and that cases
21
22 207 appear at the average rate observed in the data. That procedure was followed to overcome
23
24 208 concerns regarding the choice of uniform priors with fixed boundaries as described by Hald
25
26 209 et al. (2004), as the inferences have been shown in some cases to be sensitive to the choice of
27
28 210 boundaries (Glass et al., 2015). For the remaining hyperparameters we chose $\gamma_i = 1$, $\alpha^{(\theta)} =$
29
30 211 $\beta^{(\theta)} = 1$ to reflect weak prior knowledge.

31
32
33
34
35 212 The distribution of random effects was changed from a log-normal distribution with a mean
36
37 213 of 1 to a *Gamma* distribution with a mean of 1, which allows a Gibbs step to be used to
38
39 214 update q_i during the Markov chain Monte Carlo (MCMC) (Denison, 2002). The full
40
41 215 conditional distribution was given by

42
43
44
45 216 $q_i | (\theta, r, a, Y) \sim \text{Gamma}(\theta + \sum_{t=1}^T Y_{it}, \theta + T \sum_{j=1}^S r_{ij} a_j)$. The hyperparameter θ can be
46
47 217 fixed, but we chose to specify a prior for θ so that the variance of the random effects was
48
49 218 estimated from data.

50
51
52
53 219 The relative prevalence parameters r_{ij} were updated using Metropolis-Hastings updates
54
55 220 (Chib and Greenberg, 1995) with proposals from the prior distribution in both large and small
56
57 221 blocks. The large blocks consisted of all types $i = 1, \dots, I$ for a given source and the small
58
59
60

1
2
3 222 blocks consisted of just two types for a given source. The source dependent factors a_j were
4
5 223 updated jointly using Metropolis-Hastings random walk proposals on the log scale. To avoid
6
7 224 poor mixing when the values of a_j were small, single component Metropolis-Hastings
8
9
10 225 proposals from the prior distribution were applied. The model described above is considered
11
12 226 a single attribution model because it assumes that the source attribution is the same for each
13
14
15 227 year in the dataset.

16
17
18 228 In the temporal source attribution model, the source dependent factors depend on time and
19
20 229 therefore produce a different source attribution for each year in the study. Given that the
21
22 230 source data are sparser than the human data, and they were not collected throughout the
23
24 231 whole study period, relative prevalences r_{ij} that do not depend on time were used.

25
26
27
28 232 The temporal attribution model becomes $Y_{it} \sim Pois(\sum_{j=1}^S \lambda_{ijt})$, for year $t = 1, \dots, T$ and
29
30 233 type $i = 1, \dots, I$. The mean number of cases is decomposed as

31
32
33
34 234
$$\lambda_{ijt} = r_{ij} q_i a_{jt}.$$

35
36
37 235 Similarly, $a_{jt} \sim Exp(\alpha_{jt})$ is assumed, and $\alpha_{jt}^{-1} = \sum_{i=1}^I Y_{it}/S$, dividing the prior weight
38
39 236 equally between sources as before.

40
41
42 237 Convergence was confirmed for all parameters by visual inspection of the trace plots and
43
44 238 comparison of the posteriors from chains with randomly chosen starting values. Once
45
46
47 239 convergence was established, any evidence of lack-of-fit in the source attribution model was
48
49 240 likely due to human cases that were difficult to attribute to any of the sources in the model.
50
51 241 To help investigate that feature, an additional source, referred to as a non-sampled source,
52
53 242 was included in the model (for which no source data were observed) to identify the quantity
54
55
56 243 and profile of unattributable cases. To further assess the model performance, data from a
57
58 244 previous *Campylobacter* source attribution analysis (Mullner et al., 2009a) was analyzed and
59
60

1
2
3 245 outputs were compared.
4
5

6 246 The model was run in R (version 3.2.3; R Core Team, 2016) and to facilitate dissemination,
7
8 247 an associated R package is under development. This package containing the updated model
9
10 248 includes adaptive proposals so that unreasonable mixing is detected and the proposals are
11
12 249 adjusted to combat it. Furthermore, the package includes functions to aid interpretation of the
13
14 250 output by non-specialists.
15
16

17 18 251 Visual attribution application development 19

20
21 252 An application was created to allow users to dynamically investigate the occurrence of
22
23 253 *Salmonella* serotypes in different sources and identify patterns in the data. This application,
24
25 254 named *Source Explorer*, was developed in an RStudio Shiny (<http://shiny.rstudio.com>)
26
27 255 framework, which enabled the tool to be used locally, as a stand-alone version, or through a
28
29 256 web-based interface. The dataset of *Salmonella* serotypes per source was used as input data,
30
31 257 including the years in which the isolates were collected. Five main visualization outputs were
32
33 258 displayed, including 1) the proportion of the top 10 serotypes of each selected source over
34
35 259 time, 2) user-selected serotypes for the selected source over time, 3) a bar chart comparing
36
37 260 the top 10 serotypes of the selected source with all other sources, 4) the top 10 human
38
39 261 serotypes per year overlaid with non-human sources, and 5) rarefaction results with options
40
41 262 to select a subset of sources. The proportion of isolates per source and year was displayed to
42
43 263 account for different sampling efforts during the study period.
44
45
46
47

48
49 264
50
51
52
53
54
55
56
57
58
59
60

1
2
3 265 Results
4
5
6

7 266 *Salmonella* diversity
8

9 267 Table 2 shows the total number of *Salmonella* isolates derived from human, porcine, bovine,
10 268 turkey and chicken sources in the full dataset and the number of isolates with serotypes found
11
12 269 in both human and non-human sources. A total of 991 and 749 international travel and
13
14
15 270 outbreak associated cases were excluded, respectively. An additional 330 cases were further
16
17 271 excluded, as they were either *S. Typhi*, *S. Paratyphi*, or of an unknown serotype. The number
18
19 272 of sporadic and domestically acquired human cases per year in the dataset ranged from a
20
21 273 minimum of 405 cases in 2009 to a maximum of 546 cases in 2013. In total, 240 *Salmonella*
22
23 274 serotypes were found, 156 (65%) of which were isolated from a single source type. Twelve
24
25 275 (5%) serotypes were found in all five sources. Serotype *S. Typhimurium* var. 5- (formerly
26
27 276 known as *S. Copenhagen*) had the largest total number of isolates, with 1,617 isolates from
28
29
30 277 human, bovine, and porcine sources.
31
32
33

34
35 278 Rarefaction analysis indicated a larger serotype richness in human *Salmonella* isolates than in
36
37 279 those from non-human sources (Figure 1). In addition, a lower level of diversity was found
38
39 280 among isolates from chickens than those from other sources. Simpson's index of diversity of
40
41 281 isolates from each source is presented in Table 3. *Salmonella* isolates from human, turkey,
42
43 282 and porcine sources had the highest serotype diversity, whereas isolates from chicken and
44
45 283 bovine sources were less diverse. Overall, the PSI analysis indicated a relatively low level of
46
47 284 similarity between the different sources (Table 3). Human isolates were most similar to
48
49 285 porcine isolates, however 95% bootstrap confidence intervals overlapped between all
50
51 286 sources.
52
53
54

55
56 287
57
58
59
60

1
2
3 288 Source attribution
4

5 289 Seventy-five different serotypes were found both in humans and in at least one of the non-
6
7 290 human sources, 50 of which included at least five human cases, representing 96% of isolates
8
9
10 291 included in the full dataset. The remaining 25 serotypes included 785 isolates and were
11
12 292 designated as “other”.
13

14
15 293 The number of human salmonellosis cases attributed to each source in the single attribution
16
17 294 model is shown without (Figure 2A) and with (Figure 2B) inclusion of a non-sampled source.
18
19 295 Out of 4,672 human cases, the largest share was attributed to chickens, accounting for 2,790
20
21
22 296 (60%) and 2,049 (45%) without and with a non-sampled source, respectively.
23

24
25 297 Results from the temporal attribution model are shown in Figure 3. A marked change in the
26
27 298 attribution estimates after 2009 was observed, with an increase in the number of cases
28
29 299 attributed to the non-sampled source. The human cases most frequently assigned to the non-
30
31 300 sampled source by the model included serotypes *S.* 4,5,12:i:- (8.8%), *S.* Enteritidis (7.6%), *S.*
32
33 301 Berta (6.1%), *S.* Infantis (5.9%) and *S.* Heidelberg (4.0%).
34
35

36
37 302 The temporal model generally showed sufficient fit, aligned with previous validations of the
38
39 303 model where expected and observed cases were compared to test the validity of the model
40
41 304 outputs (e.g. Hald et al. 2004). A more stringent model assessment was also performed in
42
43 305 which the prior and the posterior distributions were compared for each of the relative
44
45 306 prevalence parameters. The prior for these parameters is based on the source data whilst the
46
47 307 posterior incorporates both the source and human data, and so if there is disagreement
48
49 308 between these two distributions then some of the human data cannot be aligned with any of
50
51 309 the sources. A selection of these plots is presented in Supporting Figure 1. Outputs from the
52
53 310 Mullner et al. (2009a) data using the current model (data not shown) were consistent with the
54
55 311 results presented in the original paper, further supporting the performance of the model
56
57
58
59
60

1
2
3 312 presented here.
4
5

6 313
7
8

9 314 Visual attribution
10

11 315 Multiple filters and visualization options facilitated extensive data exploration in *Source*

12 316 *Explorer*. Outputs were visualized in charts and data tables and were instantaneously updated

13 317 as search variables were modified. Individual serotypes or the top 10 serotypes for each

14 318 source could be selected and visualized over time. For example, the proportion of human

15 319 isolates belonging to each of the top 10 serotypes from humans was overlaid with the

16 320 proportion of isolates from each source (Figure 4). This example demonstrates the high

17 321 proportion of porcine *S. Typhimurium* var. 5- isolates throughout the study period, compared

18 322 to other sources.
19
20

21 323 *Source Explorer* enabled us to further investigate the source attribution results and potential

22 324 reasons for the spike in cases attributed to a non-sampled source after 2009. Out of the top 10

23 325 most common human serotypes in our dataset, *S. Enteritidis* and *S. 4,5,12:i:-* increased

24 326 relative to the other serotypes isolated after 2009 (Figure 5). Further, these two serotypes

25 327 were not frequently isolated from non-human sources throughout the study period,

26 328 representing 0.5 and 2.5 percent of serotypes isolated from non-human sources, respectively.
27
28
29
30

31 329
32
33
34
35
36
37
38
39
40
41
42
43
44

45 330 Discussion
46
47

48 331 Source attribution models continue to be developed and modified to accommodate different

49 332 environments, pathogens, and data sources (David et al., 2013; Glass et al., 2015; Guo et al.,

50 333 2011; Hald et al., 2004; Little et al., 2010; Mughini-Gras et al., 2014; Mullner et al., 2009a).

51 334 A detailed exploration of the data that are inputted into such models can help drive their
52
53
54
55
56
57
58
59
60

1
2
3 335 evolution to best incorporate and model observed data. Therefore, a multi-pillared approach
4
5 336 was used here to explore the burden of sporadic, domestically-acquired human *Salmonella*
6
7 337 infections in Minnesota from different sources over a ten-year period. Serotype diversity and
8
9
10 338 similarity in human, bovine, porcine, chicken, and turkey sources was assessed and a visual
11
12 339 attribution tool was developed to facilitate data exploration and validation. A modified
13
14 340 Bayesian source attribution model was subsequently used to estimate the relative contribution
15
16 341 of *Salmonella* isolates from alternative sources.
17
18
19 342 *Source Explorer* enabled users to visually assess and interact with the large datasets
20
21 343 commonly included in source attribution analyses. Attribution results here were not only
22
23 344 clarified through data exploration, but this visualization tool also helped to identify odd
24
25 345 trends in the data, such as drastic changes in the proportion of common serotypes over time.
26
27 346 For example, a lack of human *S. Typhimurim* var. 5- isolates after 2011, whereas that
28
29 347 serotype previously accounted for an average of 30% of all human *Salmonella* serotypes,
30
31 348 suggested a change in 2012. Such a finding led to further investigation into serotype naming
32
33 349 convention changes by the CDC and the state public health laboratory (CDC, 2014) and
34
35 350 ultimately led to the retrospective electronic coding of human *S. Typhimurim* var. 5- cases
36
37 351 after 2011. Quality of input data is critical to provide accurate attribution results that best
38
39 352 explain the burden of human disease. A thorough description of data and whether trends
40
41 353 accurately reflect attribution estimates is an important step that should be included in any
42
43 354 attribution analysis. That step further ensures a comprehensive and transparent assessment
44
45 355 that can be more easily communicated to various stakeholder groups to build awareness and
46
47 356 engagement and ultimately support the development of control strategies (Carroll et al.,
48
49 357 2014).
50
51
52
53
54
55
56
57
58
59
60

1
2
3 358 The highest number of sporadic and domestic cases were attributed to chickens in the single
4
5 359 attribution analyses, both with and without the inclusion of a non-sampled source. In the
6
7 360 temporal analysis, the non-sampled source had the most attributed cases, followed by
8
9 361 chickens. That finding supports previous attribution studies in the United States, where
10
11 362 poultry was the leading source of human salmonellosis (Batz et al., 2012; Gould et al., 2013;
12
13 363 Guo et al., 2011; Hoffmann et al., 2007). The temporal attribution results showed a surprising
14
15 364 amount of smoothness over time in the early part of the study, particularly considering that
16
17 365 no smoothing factor was incorporated into the model. That time period is also when most of
18
19 366 the data originated from some sources and so it may well be the case that changes in
20
21 367 attribution seen after 2010 may be the result of the source data becoming out of date rather
22
23 368 than such substantial changes in the number of cases coming from each source.
24
25
26
27

28 369 The wide credible intervals in Figure 2b show that there is a large amount of uncertainty in
29
30 370 the estimates regarding the dominant source in the single attribution model, with poultry and
31
32 371 non-sampled being the two largest estimates. In the temporal attribution model, there is a
33
34 372 period where poultry appears to make the largest contribution and a period where non-
35
36 373 sampled appears to make the largest contribution. Given that the single attribution model
37
38 374 cannot incorporate temporal changes, except as Poisson variation in the number of cases, in
39
40 375 this study, we conclude that the single attribution model does not fit as well as the temporal
41
42 376 attribution model, which results in the differences observed. This was further supported by
43
44 377 the observed vs. expected plot (Supporting Figure 2), which shows more variation in the
45
46 378 single attribution model between years than predicted by the Poisson distribution. Better
47
48 379 agreement can be observed between observed and expected cases in the temporal model.
49
50 380 Remaining deviations might indicate remaining lack of fit, but could also be explained by
51
52 381 random variation in source sampling.
53
54
55
56
57
58
59
60

1
2
3 382 As only four non-human sources had sufficient numbers of isolates to be included here, there
4
5 383 was likely an overestimation of the burden of illness from these sources. The notable
6
7 384 exclusion of eggs as a source may have skewed the attribution estimates, as eggs are a known
8
9 385 source of *Salmonella*, particularly *S. Enteritidis* (Hedberg et al., 1993; Wright et al., 2016).
10
11 386 Indeed, *S. Enteritidis* did not show a good fit between observed and posterior cases (data not
12
13 387 shown) in the current study. A European Union Food Safety Authority report ((BIOHAZ)
14
15 388 EFSA Panel on Biological Hazards 2013) recently highlighted the need to include isolates
16
17 389 from all potential major non-human sources and that the use of surrogate data or the
18
19 390 exclusion of relevant sources may seriously bias attribution results. The inclusion of a non-
20
21 391 sampled source in this study was a novel approach in source attribution analyses to deal with
22
23 392 the common problem of missing source data; however, it was previously performed as an
24
25 393 exploratory technique in other contexts and further work is necessary to confirm the validity
26
27 394 of such an approach (Pella and Masuda, 2001). The consistent results obtained using data
28
29 395 from a previous source attribution analysis (Mullner et al., 2009a) encourages the use of this
30
31 396 model in additional settings. The availability of the model in a dedicated R package will
32
33 397 facilitate this.
34
35
36
37
38

39 398 The supplementary (molecular) epidemiological analysis supported findings from the
40
41 399 attribution model. Diversity and similarity statistics highlighted a high degree of diversity and
42
43 400 low similarity of *Salmonella* serotypes isolated from all sources. Rarefaction curves for each
44
45 401 of the sources did not appear to reach a plateau, indicating that the serotype richness was not
46
47 402 fully captured by the current dataset and increased sampling effort in all sources could
48
49 403 improve the model fit. Serotype richness was greatest in humans, even when accounting for
50
51 404 sampling effort, and the PSI results indicated that overall similarity between sources was
52
53 405 relatively low. Low similarity between non-human sources can be advantageous in an
54
55 406 attribution analysis, as it supports the attribution of specific serotypes to those sources in
56
57
58
59
60

1
2
3 407 which it is uniquely found (Barco et al., 2013). In a previous study evaluating *Campylobacter*
4
5 408 spp. in New Zealand (Mullner et al., 2009b), a higher similarity among sources and humans
6
7 409 was found, with similarity estimates ranging from 0.18 to 0.58.
8
9

10 410 The most appropriate subtyping method for source attribution is one that provides an
11
12 411 appropriate level of discrimination to define subtypes associated with specific sources (EFSA
13
14 412 Panel on Biological Hazards, 2013). Serotyping is a common method for differentiation of
15
16 413 *Salmonella* types, though molecular typing methods (e.g. multiple-locus variable-number
17
18 414 tandem repeat analysis, pulsed-field gel electrophoresis) have also been used (Barco et al.,
19
20 415 2013). In some attribution studies, authors adjusted the serotype nomenclature, minimizing
21
22 416 the total number by combining variants of the same serotype. For example, Guo et al. (2011)
23
24 417 combined variants into the non-variant serotypes. This was not performed here because
25
26 418 serotype variants were assumed to not change during the course of transmission and
27
28 419 sampling. Further supporting this decision, *S. Typhimurium* and *Typhimurium* var. 5- were
29
30 420 the two most commonly isolated *Salmonella* serotypes, yet were differentially distributed
31
32 421 among sources in this study, with var 5- known to show some host association to swine
33
34 422 reservoirs; hence, that feature provided critical anchor points for the attribution model, which
35
36 423 relies on host association of subtypes. In consequence, the differentiation of the large number
37
38 424 of *S. Typhimurium* isolates into two distinct serotypes in this study improved the model fit.
39
40
41
42
43
44

45 425 Minnesota data were exclusively used in the source attribution analysis, despite much of the
46
47 426 food consumed in Minnesota likely originated from outside the state. Nevertheless, there is
48
49 427 an increasing trend in the consumption of locally produced foods, as consumers seek direct
50
51 428 farm to retail options (Low et al., 2015). Local production, import, and export data have also
52
53 429 been included in an attribution analysis to account for the flow of food commodities across
54
55 430 borders (De Knecht et al., 2015), revealing that individual countries within the European
56
57
58
59
60

1
2
3 431 Union had different attribution estimates. Further efforts in expanding data collection at a
4
5 432 state level in the United States could similarly help to elucidate spatial patterns in attribution.
6
7 433 Inclusion of data on geographical origin of isolates from retail foods, such as where the foods
8
9 434 were processed and purchased, could also potentially refine the current analysis. An
10
11 435 expanded data collection in the United States should ideally include a large number of
12
13 436 samples originating directly from food sources or food processing environments, as opposed
14
15 437 to animal reservoirs, which is where the majority of non-human isolates in this study were
16
17 438 derived. Nevertheless, given the lack of data from food sources, such data from animal
18
19 439 sources are commonly used in source attribution analyses (Mughini-Gras and van Pelt, 2014).
20
21 440 As with any analysis that relies on reported cases, underreporting of salmonellosis may have
22
23 441 introduced some bias into this study.
24
25
26
27

28 442 In summary, results here demonstrated an enhanced approach to source attribution that
29
30 443 encourages data exploration through diversity statistics and visual attribution both prior to
31
32 444 and after the use of a Bayesian source attribution model. Results here will help to inform
33
34 445 preventive and control strategies for *Salmonella* infection in Minnesota.
35
36
37

38 446
39
40

41 447 Acknowledgements

42
43
44 448 This study was funded by the Global Food Venture MnDrive initiative.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

449 References

- 450 Barco, L., Barrucci, F., Olsen, J.E., Ricci, A., 2013. *Salmonella* source attribution based on
451 microbial subtyping. *Int. J. Food. Microbiol.* 163, 193–203.
- 452 Batz, M.B., Doyle, M.P., Morris Jr, G., Painter, J., Singh, R., Tauxe, R. V, Taylor, M.R., Lo
453 Fo Wong, D.M., 2005. Attributing illness to food. *Emerg. Infect. Dis.* 11, 993–9.
- 454 Batz, M.B., Hoffmann, S., Morris Jr, J.G., 2012. Ranking the disease burden of 14 pathogens
455 in food sources in the United States using attribution data from outbreak investigations
456 and expert elicitation. *J. Food Prot.* 75, 1278–91.
- 457 Carroll, L.N., Au, A.P., Detwiler, L.T., Fu, T., Painter, I.S., Abernethy, N.F., 2014.
458 Visualization and analytics tools for infectious disease epidemiology: a systematic
459 review. *J. Biomed. Inform.* 51, 287–98.
- 460 CDC, 2014. National *Salmonella* surveillance annual report, 2012. US Dep. Heal. Hum. Serv.
461 CDC Atlanta, United States.
- 462 CDC, 2016. Foodborne Diseases Active Surveillance Network (FoodNet): FoodNet
463 Surveillance Report for 2014. CDC Atlanta, United States.
- 464 Chib, S., Greenberg, E., 1995. Understanding the metropolis-hastings algorithm. *Am. Stat.*
465 49, 327–35.
- 466 Crim, S.M., Griffin, P.M., Tauxe, R., Marder, E.P., Gilliss, D., Cronquist, A.B., Cartter, M.,
467 Tobin-D'Angelo, M., Blythe, D., Smith, K., 2015. Preliminary incidence and trends of
468 infection with pathogens transmitted commonly through food—Foodborne Diseases
469 Active Surveillance Network, 10 US sites, 2006–2014. *Morb. Mortal. Wkly. Rep.* 64,
470 495–9.
- 471 David, J.M., Guillemot, D., Bemrah, N., Thébault, A., Brisabois, A., Chemaly, M., Weill,
472 F.X., Sanders, P., Watier, L., 2013. The Bayesian microbial subtyping attribution model:
473 robustness to prior information and a proposition. *Risk Anal.* 33, 397–408.
- 474 De Knegt, L. V, Pires, S.M., Hald, T., 2015. Attributing foodborne salmonellosis in humans
475 to animal reservoirs in the European Union using a multi-country stochastic model.
476 *Epidemiol. Infect.* 143, 1175–86.
- 477 Denison, D.G.T., 2002. Bayesian methods for nonlinear classification and regression. John
478 Wiley & Sons.
- 479 Ebel, E.D., Williams, M.S., Cole, D., Travis, C.C., Klontz, K.C., Golden, N.J., Hoekstra,
480 R.M., 2016. Comparing characteristics of sporadic and outbreak-associated foodborne
481 illnesses, United States, 2004–2011. *Emerg. Infect. Dis.* 22, 1193.
- 482 EFSA Panel on Biological Hazards, 2013. Scientific Opinion on the evaluation of molecular
483 typing methods for major food-borne microbiological hazards and their use for
484 attribution modelling, outbreak investigation and scanning surveillance: Part 1
485 (evaluation of methods and applications). *EFSA J.* 11, 3502.
- 486 Glass, K., Fearnley, E., Hocking, H., Raupach, J., Veitch, M., Ford, L., Kirk, M.D., 2015.

- 1
2
3 487 Bayesian source attribution of salmonellosis in South Australia. *Risk Anal.* 36, 561-70.
4
5 488 Gould, L.H., Walsh, K.A., Vieira, A.R., Herman, K., Williams, I.T., Hall, A.J., Cole, D.,
6 489 2013. Surveillance for foodborne disease outbreaks—United States, 1998–2008.
7 490 *MMWR Surveill. Summ.* 62, 1–34.
8
9 491 Grimont, P.A.D. and Weill, F.X., 2007. Antigenic formulae of the *Salmonella* serovars.
10 492 WHO collaborating centre for reference and research on *Salmonella*, 9.
11
12 493 Guo, C., Hoekstra, R.M., Schroeder, C.M., Pires, S.M., Ong, K.L., Hartnett, E., Naugle, A.,
13 494 Harman, J., Bennett, P., Cieslak, P., 2011. Application of Bayesian techniques to model
14 495 the burden of human salmonellosis attributable to US food commodities at the point of
15 496 processing: adaptation of a Danish model. *Foodborne Pathog. Dis.* 8, 509–16.
16
17
18 497 Hald, T., Vose, D., Wegener, H.C., Koupeev, T., 2004. A Bayesian approach to quantify the
19 498 contribution of animal-food sources to human salmonellosis. *Risk Anal.* 24, 255–69.
20
21 499 Hammer, Ø., Harper, D.A.T., Ryan, P.D., 2009. PAST-PAleontological STatistics, ver.
22 500 1.89. University of Oslo, Oslo 1–31.
23
24 501 Hedberg, C.W., David, M.J., White, K.E., MacDonald, K.L., Osterholm, M.T., 1993. Role of
25 502 egg consumption in sporadic *Salmonella* enteritidis and *Salmonella* typhimurium
26 503 infections in Minnesota. *J. Infect. Dis.* 167, 107–11.
27
28
29 504 Hoffmann, S., Fischbeck, P., Krupnick, A., McWilliams, M., 2007. Using expert elicitation to
30 505 link foodborne illnesses in the United States to foods. *J. Food Prot.* 70, 1220–9.
31
32 506 Hong, S., Rovira, A., Davies, P., Ahlstrom, C., Muellner, P., Rendahl, A., Olsen, K., Bender,
33 507 J.B., Wells, S., Perez, A., Alvarez, J., 2016. Serotypes and Antimicrobial Resistance in
34 508 *Salmonella enterica* Recovered from Clinical Samples from Cattle and Swine in
35 509 Minnesota, 2006 to 2015. *PloS ONE*, 11(12), p.e0168016.
36
37 510 Kingsley, R.A., Bäumlner, A.J., 2000. Host adaptation and the emergence of infectious
38 511 disease: the *Salmonella* paradigm. *Mol. Microbiol.* 36, 1006–14.
39
40 512 Kirk, M.D., Pires, S.M., Black, R.E., Caipo, M., Crump, J.A., Devleeschauwer, B., Döpfer,
41 513 D., Fazil, A., Fischer-Walker, C.L., Hald, T., 2015. World Health Organization
42 514 estimates of the global and regional disease burden of 22 foodborne bacterial, protozoal,
43 515 and viral diseases, 2010: a data synthesis. *PLoS Med* 12, e1001921.
44
45
46 516 Little, C.L., Pires, S.M., Gillespie, I.A., Grant, K., Nichols, G.L., 2010. Attribution of human
47 517 *Listeria monocytogenes* infections in England and Wales to ready-to-eat food sources
48 518 placed on the market: adaptation of the Hald *Salmonella* source attribution model.
49 519 *Foodborne Pathog. Dis.* 7, 749–56.
50
51 520 Low, S.A., Adalja, A., Beaulieu, E., Key, N., Martinez, S., Melton, A., Perez, A., Ralston, K.,
52 521 Stewart, H., Suttles, S., Vogel, S., Jablonski, B.B.R. 2015. Trends in US local and
53 522 regional food systems: A report to Congress.
54
55
56 523 Minnesota Department of Health, 2015. Annual Summary of Communicable Disease
57 524 Reported to the Minnesota Department of Health, 2015. *Disease Control Newsletter*, 43.
58
59
60

- 1
2
3 525 Muellner, P., Marshall, J.C., Spencer, S.E.F., Noble, A.D., Shadbolt, T., Collins-Emerson,
4 526 J.M., Midwinter, A.C., Carter, P.E., Pirie, R., Wilson, D.J., 2011. Utilizing a
5 527 combination of molecular and spatial tools to assess the effect of a public health
6 528 intervention. *Prev. Vet. Med.* 102, 242–53.
- 8 529 Mughini-Gras, L., Barrucci, F., Smid, J.H., Graziani, C., Luzzi, I., Ricci, A., Barco, L.,
9 530 Rosmini, R., Havelaar, A.H., Van Pelt, W., 2014. Attribution of human *Salmonella*
10 531 infections to animal and food sources in Italy (2002–2010): adaptations of the Dutch and
11 532 modified Hald source attribution models. *Epidemiol. Infect.* 142, 1070–82.
- 13 533 Mughini-Gras, L. and van Pelt, W., 2014. *Salmonella* source attribution based on microbial
14 534 subtyping: Does including data on food consumption matter? *Int. J. Food Microbiol.*
15 535 191, 109–15.
- 17 536 Mullner, P., Jones, G., Noble, A., Spencer, S.E.F., Hathaway, S., French, N.P., 2009a. Source
18 537 attribution of food-borne zoonoses in New Zealand: a modified Hald model. *Risk Anal.*
19 538 29, 970–84.
- 21 539 Mullner, P., Spencer, S.E.F., Wilson, D.J., Jones, G., Noble, A.D., Midwinter, A.C., Collins-
22 540 Emerson, J.M., Carter, P., Hathaway, S., French, N.P., 2009b. Assigning the source of
23 541 human campylobacteriosis in New Zealand: a comparative genetic and epidemiological
24 542 approach. *Infect. Genet. Evol.* 9, 1311–9.
- 26 543 Painter, J.A., Hoekstra, R.M., Ayers, T., Tauxe, R. V, Braden, C.R., Angulo, F.J., Griffin,
27 544 P.M., 2013. Attribution of foodborne illnesses, hospitalizations, and deaths to food
28 545 commodities by using outbreak data, United States, 1998–2008. *Emerg. Infect. Dis.* 19,
29 546 407–15.
- 31 547 Pella, J. and Masuda, M., 2001. Bayesian methods for analysis of stock mixtures from
32 548 genetic characters. *Fishery Bulletin*, 99(1), pp.151-167.
- 34 549 Pires, S.M. and Hald, T., 2010. Assessing the differences in public health impact of
35 550 *Salmonella* subtypes using a Bayesian microbial subtyping approach for source
36 551 attribution. *Foodborne Pathog. Dis.* 7, 143–51.
- 38 552 R Core Team, 2016. R: A language and environment for statistical computing. R Foundation
39 553 for Statistical Computing, Vienna, Austria.
- 41 554 Scallan, E., Hoekstra, R.M., Angulo, F.J., Tauxe, R. V, Widdowson, M.-A., Roy, S.L., Jones,
42 555 J.L., Griffin, P.M., 2011. Foodborne illness acquired in the United States—major
43 556 pathogens. *Emerg. Infect. Dis.* 17.
- 45 557 Sears, A., Baker, M.G., Wilson, N., Marshall, J., Muellner, P., Campbell, D.M., Lake, R.J.,
46 558 French, N.P., 2011. Marked campylobacteriosis decline after interventions aimed at
47 559 poultry, New Zealand. *Emerg. Infect. Dis.* 17, 1007–15.
- 49 560 Uzzau, S., Brown, D.J., Wallis, T., Rubino, S., Leori, G., Bernard, S., Casadesús, J., Platt,
50 561 D.J., Olsen, J.E., 2000. Host adapted serotypes of *Salmonella enterica*. *Epidemiol.*
51 562 *Infect.* 125, 229–55.
- 53 563 Wright, A.P., Richardson, L., Mahon, B.E., Rothenberg, R., Cole, D.J., 2016. The rise and
54 564 decline in *Salmonella enterica* serovar Enteritidis outbreaks attributed to egg-containing

1
2
3 565 foods in the United States, 1973-2009. *Epidemiol. Infect.* 144, 810.
4
5 566
6

7 567 Supporting Information 8 9

10
11 568 Supporting Figure 1. Posterior distribution of relative prevalence of *S. Enteritidis* in chicken
12
13 569 with informative prior density (based on source typing data) shown in red for four different
14
15 570 models – A) Single attribution including only sampled sources, B) Single attribution
16
17 571 including non-sampled sources, C) Temporal attribution including only sampled sources, and
18
19 572 D) Temporal attribution including non-sampled sources.
20
21
22

23 573
24
25

26 574 Supporting Figure 2 Plots of the observed and expected cases for individual *Salmonella*
27
28 575 serotypes in the A) Single attribution model including only sampled sources, B) Single
29
30 576 attribution including non-sampled sources, C) Temporal attribution including only sampled
31
32 577 sources, and D) Temporal attribution including non-sampled sources.
33
34
35

36 578
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. Parameter interpretations for the source attribution model.

Parameter	Parameter Interpretation
λ_{ij}	Expected number of human cases of type i from source j .
q_i	Strain-specific factor for strain i (e.g. survivability, pathogenicity to humans, virulence).
a_j	Source-specific factor for source j (e.g. risk through typical storage and preparation of food from source j).
r_{ij}	The relative prevalence of type i in source j .

Table 2. Total number of *Salmonella* isolates recovered from different sources and the number of isolates belonging to a serotype found in both human and food/animal sources in Minnesota.

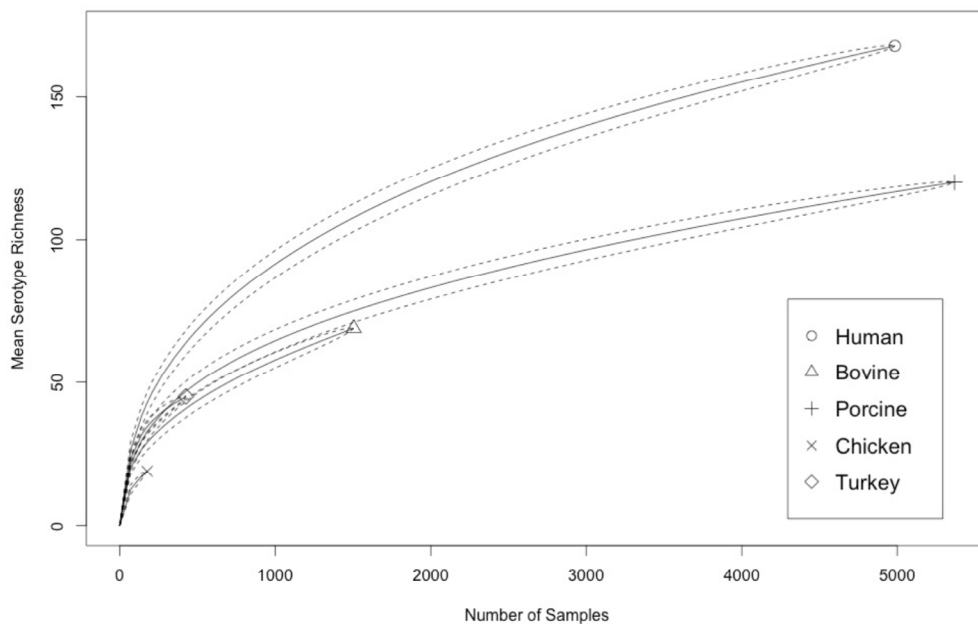
Source	# isolates in the full dataset	# isolates with shared serotypes (% of full dataset)
Human	4985	4672 (94%) ^a
Porcine	5368	5257 (98%) ^b
Bovine	1505	1484 (99%) ^b
Turkey	426	402 (94%) ^b
Chicken	177	170 (96%) ^b

^a shared with any other source

^b shared with human isolates

Table 3. Proportional similarity index and Simpson's index of diversity with 95% bootstrap CI given in parentheses of *Salmonella* serotypes in Minnesota.

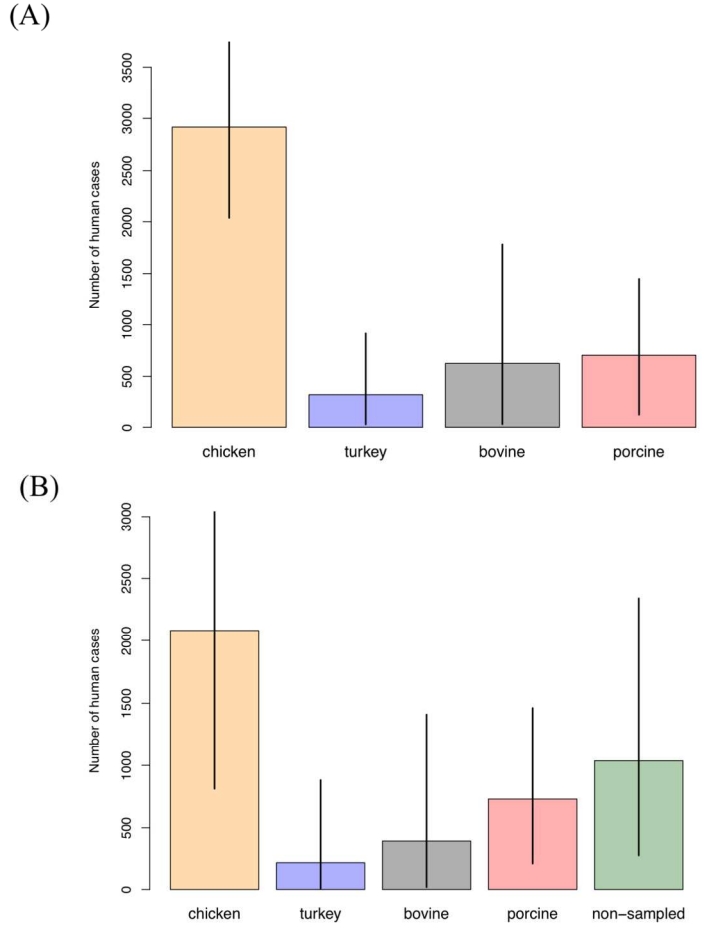
Source	PSI				Simpson's index of diversity
	Human	Bovine	Porcine	Chicken	
Human	-				0.92 (0.92, 0.93)
Bovine	0.28 (0.26, 0.30)	-			0.88 (0.87, 0.90)
Porcine	0.36 (0.34, 0.37)	0.23 (0.21, 0.25)	-		0.90 (0.90, 0.91)
Chicken	0.32 (0.24, 0.37)	0.16 (0.12, 0.19)	0.18 (0.15, 0.20)	-	0.73 (0.68, 0.79)
Turkey	0.21 (0.19, 0.24)	0.19 (0.16, 0.22)	0.30 (0.26, 0.32)	0.24 (0.18, 0.27)	0.92 (0.91, 0.93)



Rarefaction curve indicating the mean serotype richness of Salmonella serotypes from human, bovine, porcine, chicken, and turkey sources in Minnesota.

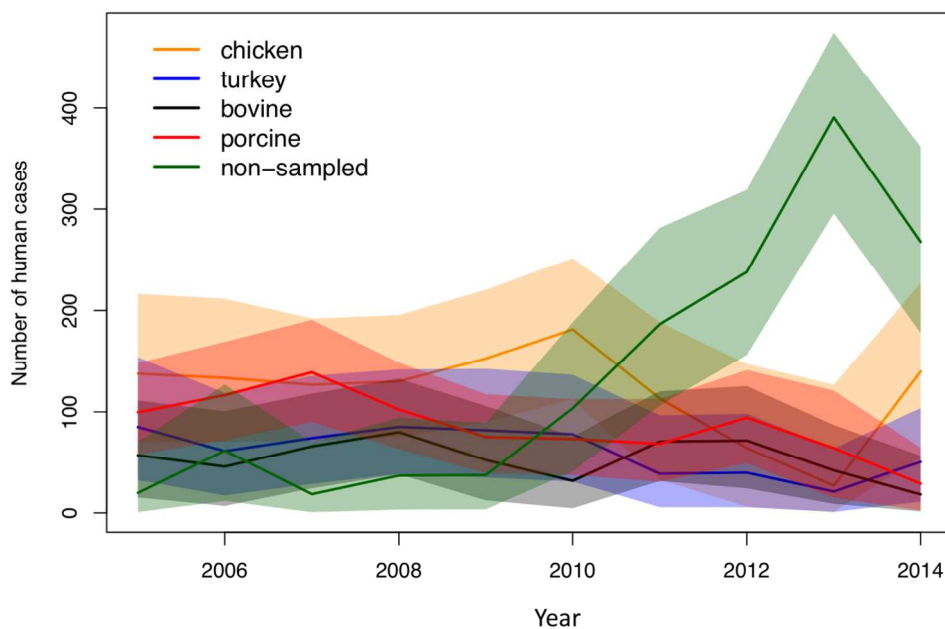
Figure 1
581x401mm (72 x 72 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Single attribution results based on 4,672 human salmonellosis cases in Minnesota between 2005 and 2014. The graphs show the number of attributed cases to each source with 95% Bayesian credible intervals without (A) and with (B) the inclusion of a non-sampled source. Year of isolation was not included in this model.

Figure 2
594x792mm (72 x 72 DPI)

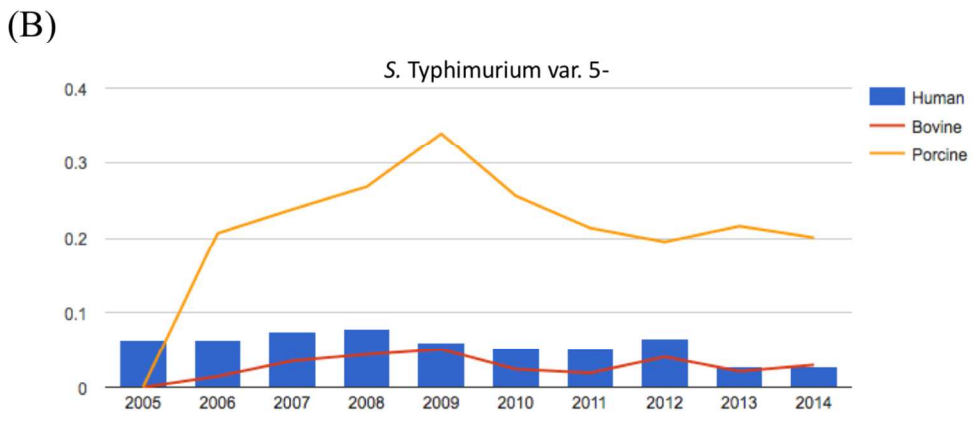
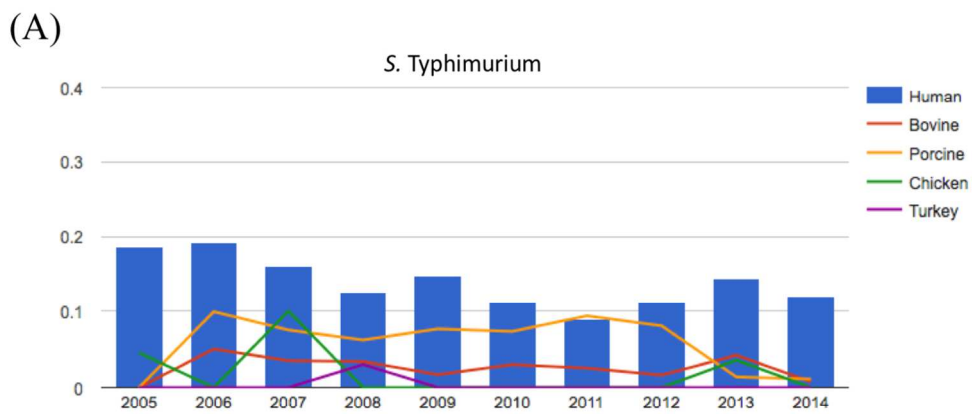


Temporal source attribution results based on 4,672 salmonellosis cases in Minnesota between 2005 and 2014. The graphs show the number of attributed cases to each source with 95% Bayesian credible intervals, incorporating the year in which human isolates were derived.

Figure 3

538x358mm (72 x 72 DPI)

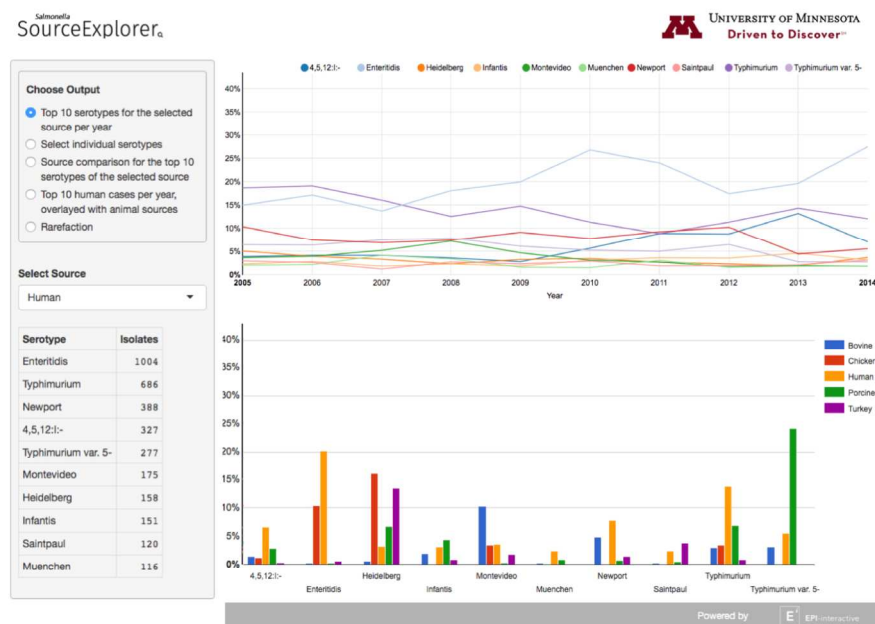
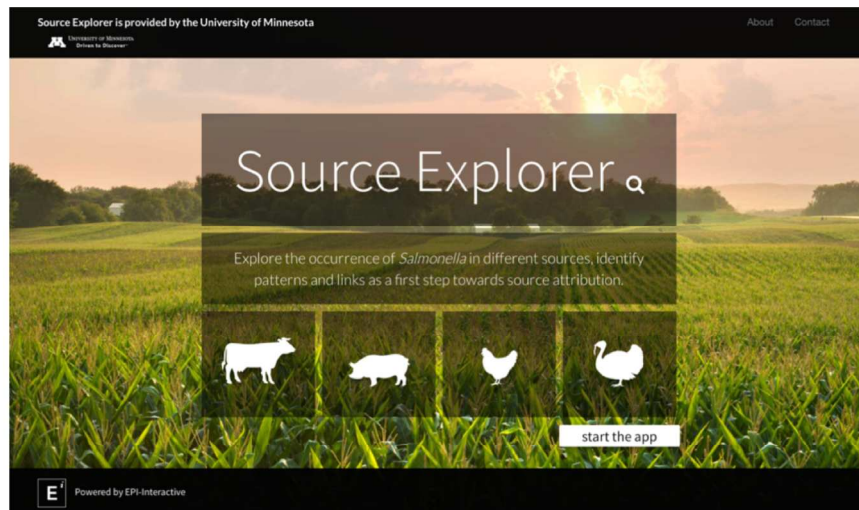
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Source Explorer outputs displaying the proportion of *S. Typhimurium* (A) and *S. Typhimurium* var. 5- (B) isolates from human and non-human sources in Minnesota between 2005 and 2014.

Figure 4
528x483mm (72 x 72 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Top panel: Source Explorer landing page. Bottom panel: Screenshot of selected Source Explorer functionalities used to explore source attribution input data and results.

Figure 5
594x792mm (72 x 72 DPI)