

A federated learning approach to data sharing in a food supply chain

Anand Gavai (✉ anand.gavai@wur.nl)

Wageningen Food Safety Research (WFSR) <https://orcid.org/0000-0002-4738-190X>

Yamine Bouzembrak (✉ yamine.bouzembrak@wur.nl)

Wageningen Food Safety Research (WFSR) <https://orcid.org/0000-0001-8028-0847>

Wenjuan Mu (✉ wenjuan.mu@wur.nl)

Wageningen University and Research (WUR)

Frank Martin (✉ F.Martin@iknl.nl)

Netherlands Comprehensive Cancer Organization

Rajaram Kaliyaperumal (✉ r.kaliyaperumal@lumc.nl)

Leiden University Medical Center

Johan van Soest (✉ johan.vansoest@maastro.nl)

Maastricht University Medical Centre

Ananya Choudhury (✉ Ananya.Choudhury@maastro.nl)

Maastricht University Medical Centre

Jaap Heringa (✉ j.heringa@vu.nl)

VU University Amsterdam

Andre Dekker (✉ andre.dekker@maastro.nl)

Maastricht University Medical Centre

Hans Marvin (✉ hans.marvin@wur.nl)

Wageningen University and Research (WUR)

Article

Keywords:

DOI: <https://doi.org/10.21203/rs.3.rs-2350301/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

A federated learning approach to data sharing in a food supply chain

Authors : Anand Gavai¹, Yamine Bouzembrak¹, Wenjuan Mu¹, Frank Martin², Rajaram Kaliyaperumal⁶, Johan van Soest^{3,4}, Ananya Choudhury⁴, Jaap Heringa⁵, Andre Dekker⁴ & Hans J.P. Marvin¹

Affiliations:

1. Wageningen Food Safety Research, Akkermaalsbos 2, 6708, WB, Wageningen, the Netherlands
2. Netherlands Comprehensive Cancer Organization (IKNL), Eindhoven, NL.
3. Brightlands Institute for Smart Society, Faculty of Science and Engineering, Maastricht University, Heerlen, The Netherlands
4. Department of Radiation Oncology (Maastr), GROW School for Oncology and Reproduction, Maastricht University Medical Centre, Maastricht, NL.
5. Centre for Integrative Bioinformatics (IBIVU), VU University Amsterdam, Amsterdam, The Netherlands.
6. Department of Human Genetics, Leiden University Medical Center, Leiden, The Netherlands.

Abstract

Ensuring safe and healthy food is a big challenge due to the complexity of food supply chains and their vulnerability to many internal and external factors. Recent research has shown that Artificial Intelligence (AI) based algorithms, in particularly data driven Bayesian Network (BN) models, are very suitable as a tool to predict future food safety risks and hence allowing food producers to take proper actions to avoid food safety problems to occur. Such models become even more powerful when data can be used from all actors in the supply chain (e.g., farmers, food producers, authorities) but data sharing is hampered by different interests, data security & data privacy. Federated learning (FL) may circumvent these problems as demonstrated in various areas of the life sciences mainly using linear models.

In this research, a federated BN was developed for food fraud to demonstrate the potential of FL for the whole food safety domain. This concept consisted of three geographically different data stations hosting different sets of food fraud data which have been made FAIR (e.g., Findable, Accessible, Interoperable & Reusable). It is demonstrated that a BN model can be

34 trained on the data of different data stations while the data never leaves its data station abiding
35 security and sensitivity requirements and that this BN model performs like the BN model
36 trained on the complete data set pulled into one data station. We demonstrated for the first time
37 the applicability of the federated BN in food safety and anticipate that such concept may support
38 stakeholders in the food supply chain for better decision-making regarding food safety and food
39 fraud control.

40

41 **Keywords:** Big data, Data Analytics, Semantic Interoperability, Bayesian Network, Federated
42 Learning, FAIR Data, Food Safety

43

45 **Introduction**

46 To ensure safe and healthy food, all actors in food supply chains are collecting large amounts
47 of food safety and quality data at the various production stages. Implementation of new
48 technologies such as drones (Almaki, 2020), mobile devices (Nelis et al. 2020) and Internet of
49 things (IoT) (Bouzembrak et al. 2019) urge the implementation of Big Data solutions in food
50 production including food safety (Marvin et al. 2017; Jin et al 2020). Artificial Intelligence (AI)
51 will play a key role in the digital transformation of the food supply chains in particularly in the
52 exploitation of these vast data sources and to support the implementation of a holistic approach
53 to ensure truly sustainable food systems (Marvin et al 2022). It was demonstrated for food
54 safety and food fraud that the AI method Bayesian Network (BN) is suitable to implement the
55 holistic approach in which data from different origins and nature are integrated (Marvin et al
56 2016, Marvin & Bouzembrak 2020, Wang et al 2022).

57 However, the impact of these new data-generating technologies depends critically on data
58 sharing and integration, which is one of the biggest challenges within the food supply chain
59 (Top et al. 2022). Different data owners have different interests and priorities that hinder the
60 incentive to share data. Data collected in the context of food safety can be politically sensitive
61 and considered a competitive advantage (Curry, 2016), but there is also a cost associated with
62 collecting this data. Nonetheless, there is a strong shared interest among stakeholders in food
63 safety compliance. It is also desirable that clear guidelines for data sharing be agreed upon. To
64 this end, extensive negotiation among data owners is usually required to resolve issues of
65 ownership, confidentiality, and management of the data. Agreement is particularly difficult
66 when many supply chain actors with conflicting interests are involved (e.g., competitors,
67 control authorities, and manufacturers, etc.). An additional challenge in data sharing is that data

68 must be described with metadata using ontologies so that it can be found by specific search
69 engines. However, aside from FOONON (Dooley et al., 2018), few food safety ontologies are
70 publicly available to date, making the adoption of data sharing and integration technologies
71 difficult. Evidence of solving these issues can be found in literature using federated learning.
72 In a federated environment, data never leave the physical location of the data owners. Instead,
73 the algorithm moves between these locations (i.e., data stations) and collects parameters from
74 the data at the data station's physical location. One of the main advantages of this approach is
75 that the federated infrastructure can perform some of the "negotiation" (otherwise done by
76 humans) automatically once data sharing policies are agreed upon. Federated learning has
77 recently gained attention in several domains such as life sciences (Flores et al., 2021; Dayan et
78 al., 2021; Beyan et al., 2020; Deist et al., 2020; Geleijnse et al., 2020; Johan van Soest, 2018;
79 Shi et al., 2019) but has, to our knowledge, not yet been explored in the food safety domain.
80 One of the key requirements is that the data at the different sites partially or fully conform to
81 the principles for discoverable, accessible, interoperable, and reusable (FAIR) data (Wilkinson
82 et al., 2016).

83 The FAIR data principles provide guidelines for individual data station owners to make their
84 internal data FAIR, which enables machines to automatically find, access, interact with, and
85 reuse data without human intervention (Top et al 2022).

86 In this study, the federated learning concept was developed to predict food fraud type through
87 federated food fraud data stations and a BN model. It was shown that a BN model could be
88 trained on these data stations without the data leaving the data stations and that the model
89 performance is similar to a BN model developed on the same data pooled to one location. The
90 developed FL infrastructure addresses some of the limitations that classical centralized
91 solutions still faces such as data ownership, confidentiality, privacy, security, and increased
92 data traffic by: (i) keeping the food data locally with the data owner; (ii) requiring no exchange

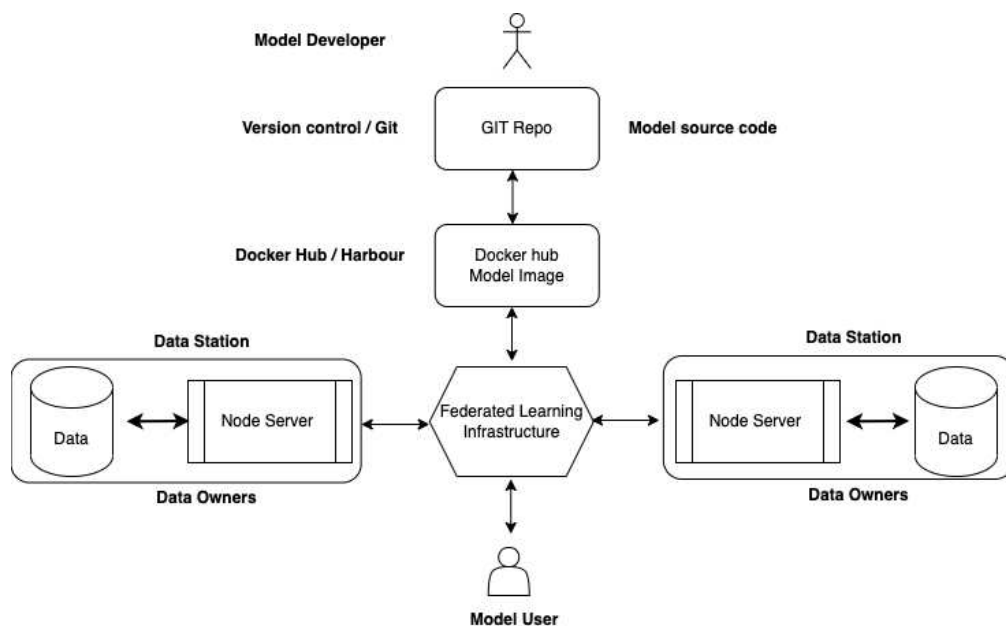
93 of raw data (iii) providing high-level data security; (iv) reducing data traffic between actors in
94 the food chain; and (v) allowing parameter learning from all stations.

95 We advocate that such approach may help to maintain the high standards needed in the food
96 production regarding food safety and food quality by exploiting all available data that is
97 distributed over the various stakeholders. The individual elements of this framework are
98 described in the Materials and Methods section of this article.

99

100 **Materials and Methods**

101 To demonstrate the concept, a minimal set of technical requirements was implemented, as
102 shown in Figures 1 & 2 (Moncada-Torres et al., 2020), which are described in more detail in
103 the next section.



104
105

106 **Fig.1:** FAIR federate architecture consisting of the following main components: model
107 developers & users, federate learning infrastructure and two data stations (i.e., data base of the
108 data owner and compute node depicted as node server).

109

110 **Federated architecture**

111 In this study, the Vantage6¹ platform (version 2.3.4) was used, which is a federated learning
112 infrastructure for secure information sharing. This central infrastructure component
113 (authentication & message broker) was hosted at the Wageningen University & Research
114 premise. Vantage6 enforces privacy concerns by allowing only certain algorithms to run. This
115 ensures that data is secure even if the security of the server is compromised. Collaboration
116 policies for data sharing were defined on the central server. To set up a federated collaboration
117 (shown in Fig.1), the following activities should be carried out: i) A collaboration network is
118 created when all participating organizations agree to work together on a particular issue. These
119 organizations can represent any actor in the food supply chain (e.g., farmers, food industry,
120 government agencies), who can be both owners and/or users of data and models/algorithms.
121 This infrastructure is created from a central location where a collaboration server is established
122 that has an integrated database that stores collaboration information and policies that determine
123 which collaborations have access to which data stations according to those policies. An
124 administrator makes individual organizations part of this collaboration at the central server and
125 distributes the authentication information to all involved organizations and users; ii) A data
126 analysis algorithm or model learning application is created by a model developer using an
127 appropriate language (e.g., R/Python). These scripts can be used on a particular data station that
128 is part of this collaboration network. These scripts are typically published as Docker images in
129 an internal Docker registry or a publicly accessible Docker registry, approved within the
130 collaboration. All input parameters are passed to this Docker image; iii) Any data station
131 requesting this Docker image as part of a collaboration can have it run at the data station owner's
132 site if execution of that image is allowed; iv) The computing nodes of the data stations return

¹ <https://vantage6.ai/>

133 the results after the algorithm has been executed. These results are sent to the central server,
134 from which the model users can retrieve the results.

135

136 **Data station architecture**

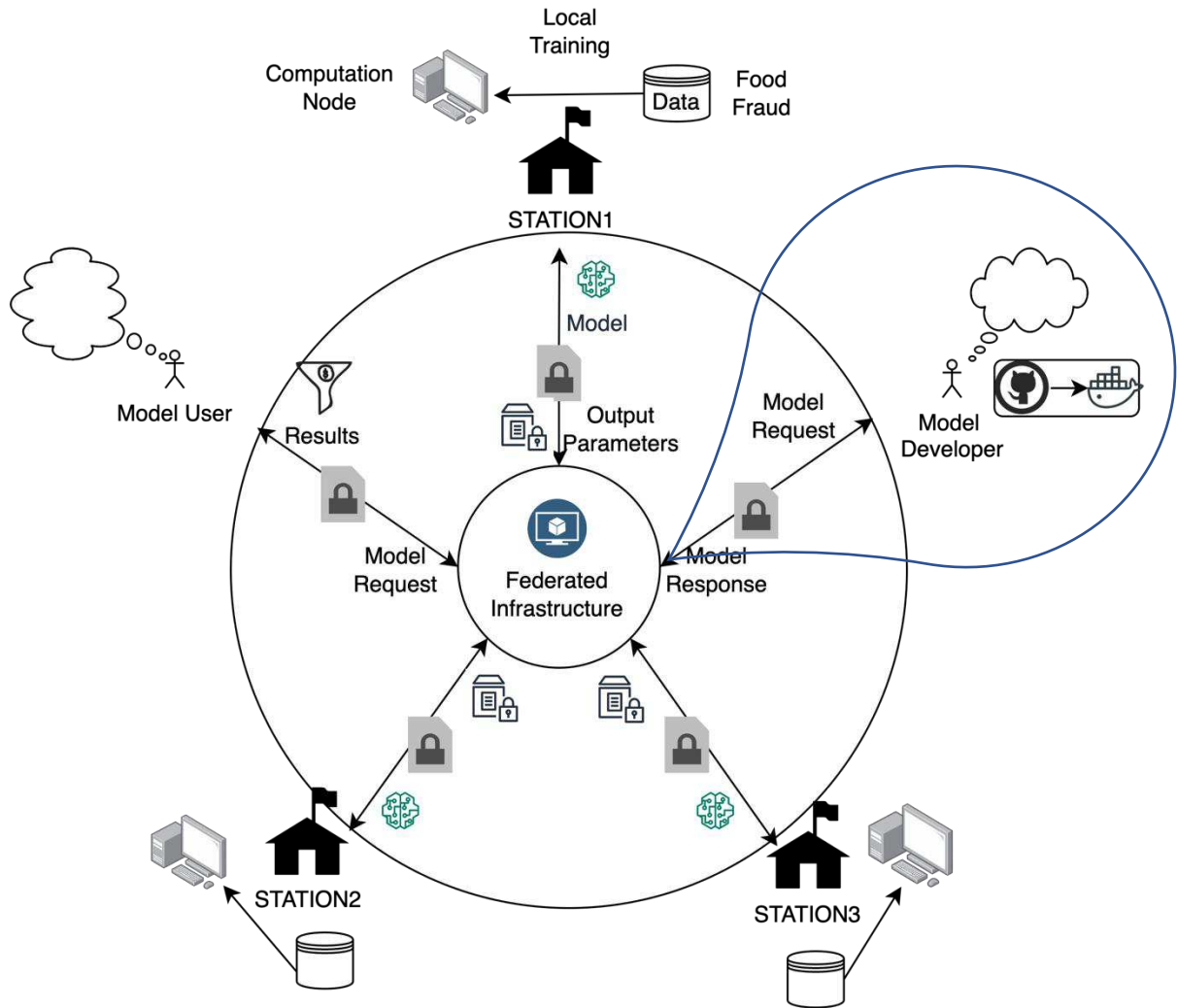
137 Each data station is connected to various components that enable the entire workflow of data
138 storage and computation in a secure environment. These components include the data volume
139 (i.e., local storage for data in csv files or a database), the Docker daemon (i.e., software to run
140 Docker images), the algorithm container (i.e., compute node running the Docker image of the
141 algorithm), CLI (i.e., command line interface to start, stop and debug nodes using log files or
142 configuring new nodes), and a set of configuration files consisting of data policies and API keys
143 that connect this data station to the federated server it belongs.

144

145 **Description of data at each data stations**

146 Data for this study is derived from two major sources: the European Union (EU) Rapid Alert
147 for Food and Feed (RASFF)² database and the United States (US) Economic Motivation
148 Adulteration (EMA) database from the period these were publicly available. The collected data
149 was separated in 3 portions and placed in different databases, representing a hypothetical
150 situation of three data owners. For each dataset, a data station was prepared and hosted at a
151 different geographic location in the Netherlands, namely Wageningen (STATION-1),
152 Maastricht (STATION-2), and Utrecht (STATION-3) (Fig 2). STATION-1 contained RASFF
153 data from 2008 to 2013 (i.e., 202 observations), STATION-2 contained RASFF data from 2014
154 to 2018 (i.e., 144 observations), while STATION-3 contained EMA data from 2008 to 2017
155 (i.e., 95 observations).

² <https://webgate.ec.europa.eu/rasff-window/screen/search>



156
 157 **Fig.2.** Conceptual framework of collaboration aimed at prediction of food fraud. The model
 158 users first train their model on their own local data stations using a Bayesian Network (BN)
 159 model (embedded inside a docker image). The trained model parameters are then transferred to
 160 the central server. These parameters are subsequently retrieved back by model users in a secured
 161 setting via the central server for each station to generate a combined BN model which contains
 162 information from each data station.

163
 164 For each of these datasets, food fraud type, product category, year, origin country and the
 165 control country were selected to be used for data training (i.e., the BN model). The meanings,
 166 corresponding nodes and states of the variables of the BN model are listed in Table 1.

167
 168

169 **Table 1. Metadata for data stations**

Variable name	Node name	States
Type of fraud	Fraud type	Artificial enhancement/Improvement, Smuggling-Mislabeling-Origin Masking, Substitution-Dilution
Category of product	Product	Alcoholic, Fish_Seafood, etc.
Year	Year	2008, 2009, ..., 2018
Origin country	Country (O)	South Korea, Croatia, etc.
Report country	Country (N)	United Kingdom, Netherlands, etc.

170

171 **Data formats**

172 All metadata belonging to each of these stations is made available on an internally hosted FAIR
 173 Data Point (da Silva Santos, et al., 2022) that can be accessed at WUR³. To demonstrate the
 174 applicability of the infrastructure, different data formats were used. The data located at
 175 STATION-1 and STATION-2, containing RASFF food fraud notifications, were in CSV format
 176 and the data located at STATION-3, containing EMA food fraud notifications, was in RDF
 177 format. Because the BN model only can consume data from CSV, the data on this station is
 178 automatically converted into a CSV format before the BN model is trained. In the situation
 179 where all data stations would host data in RDF format.

180

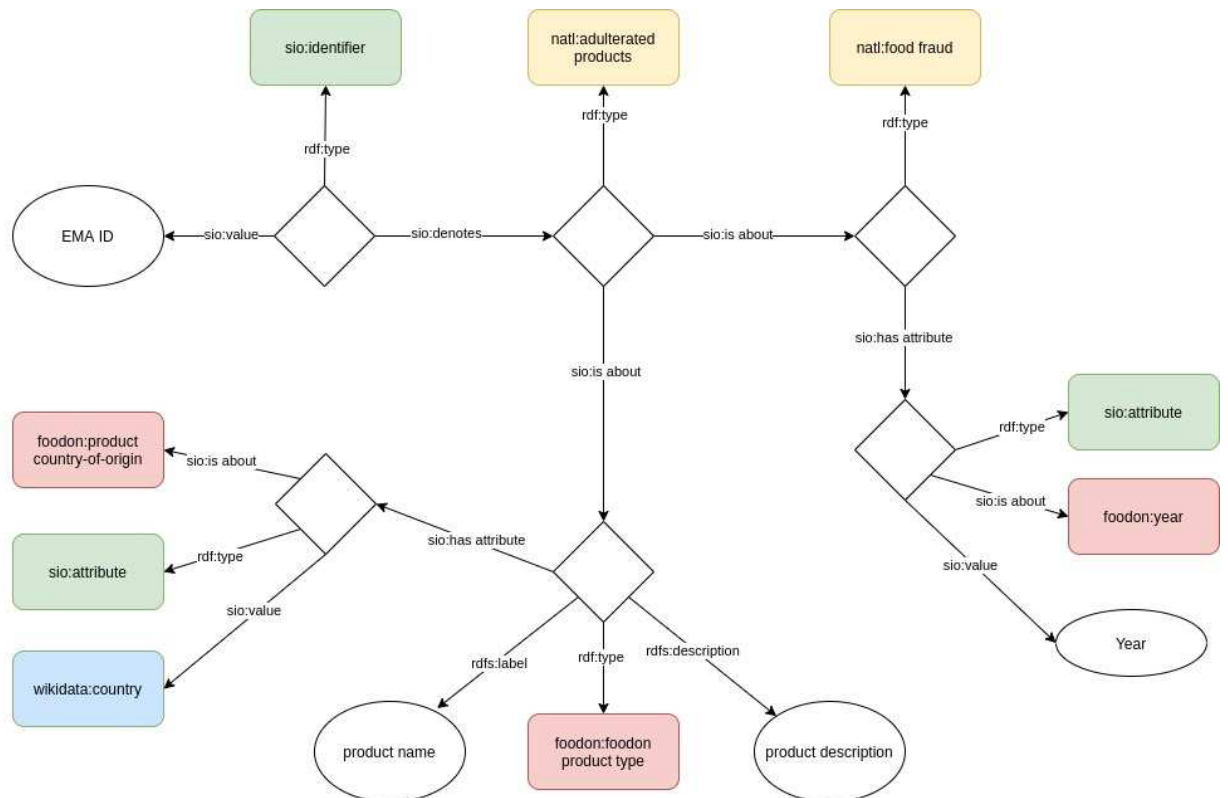
181 **Semantic interoperability**

182 The data from STATION-3 were in a semantically interoperable RDF format and are shown in
 183 Fig.3, including their metadata. This data model was created using ontologies such as the
 184 National Agricultural Library Thesaurus, FOODON and Semantic Science Integrated Ontology
 185 (SIO), and Wikidata. To date, there is no specific ontology for food safety and only minimal

³ <https://fdp.containers.wurnet.nl/>

186 information standards exist that have not yet been formalized. Therefore, we have relied on
 187 generic terminologies that would increase the use(ability) of them for the community. The
 188 AgroPortal⁴ was used to search for these ontologies, which contains various ontologies in the
 189 field of agricultural and plant sciences. This RDF data model was then converted to a CSV
 190 format by STATION-3 to be hosted on STATION-3 nodes.

191



- FoodOn ontology
- Semanticscience Integrated Ontology
- National Agricultural Library Thesaurus
- Wikidata entries
- RDF instance
- Ontology classes or other instance IRI
- RDF literal

192
193

194 **Fig.3.** RDF model of EMA data as hosted on STATION-3 node with associated ontologies
 195 with variables Year, Product, Notification, Origin, and Food Fraud type.

196

⁴ <http://agroportal.lirmm.fr/>

197

198 **Description of the BN model**

199 In this work, we implemented a BN model, in a federated environment. The R package
200 "bnlearn" (Scutari, 2010) was used to apply "Bayesian network analysis" to the data. The
201 algorithm is encapsulated in a Docker image. In this study, the R algorithm library of vantage6
202 is used. This library contains auxiliary functions for input/output between the infrastructure and
203 the algorithm. This allows developers to focus on implementing the algorithm and worry less
204 about the infrastructure-specific code. The source code for this library can be found in the
205 repository (see supplement).

206 A BN model was developed that automatically divides the dataset into a training (80%) and a
207 test dataset (20%) and learns using a data framework. The algorithms use a standard CSV file
208 as input. Tree-Augmented Naive Bayes (TAN) was applied to learn the structure of BN on each
209 of the training datasets for the variable "Fraud" (i.e., the fraud type). Once the structure was
210 known, the parameters were estimated using the bn.fit function with the "Bayes" method to
211 derive the three BN models. The three derived BN models were applied to the corresponding
212 test data sets to make predictions about "fraud". The combined ROC was calculated based on
213 the micro-average curve ROC because the fraud types in each data station were highly
214 imbalanced. Finally, the predicted "fraud" type was compared to the observed "fraud" type
215 recorded in the test data sets to obtain the prediction accuracy. This model provides the results
216 in web-enabled json format. The model was trained using characteristics such as the product
217 susceptible to fraud (eggs, oils, oily fish, and seafood), the year (when it was first reported), the
218 origin (which country the food came from), and the type of fraud (e.g., substitution, mislabeling,
219 etc.). All data stations that provided this data had these variables in common.

220 To demonstrate the operation of the developed federated infrastructure, two experiments were
221 conducted to (a) first, test how a BN model trained on the aggregate dataset of all three data

222 stations in a federated environment performs on the test dataset of each of the data stations
 223 separately. Second, (b) we tested how a model from BN, trained on the aggregated dataset of
 224 all three data stations without a federated environment, performs on the aggregated test dataset.

225

226 **Results and Discussion**

227 **Food fraud data**

228 Often actors in the supply chain have limited, imbalanced data available on which decisions
 229 must be made. Sharing these datasets between these actors would improve their individual
 230 models and decision making. To mimic such situation and to demonstrate how FL may solve
 231 this data sharing issue, the total available data set was separated into three incomplete data sets
 232 varying in the number of food fraud cases, years and type of fraud (see Table 2 and Fig.4). For
 233 example, STATION-1 contained only two types of food fraud, which are smuggling-
 234 mislabeling-origin masking (i.e., 105 observations), and substitution-dilution (i.e., 97
 235 observations). All the data of STATION-1 is obtained from RASFF from 2008 to 2013 (see
 236 Table 2). The other stations contained other types of food fraud (i.e., Artificial
 237 enhancement/improvement) and more recent food fraud data (e.g., 2014-2018).

238

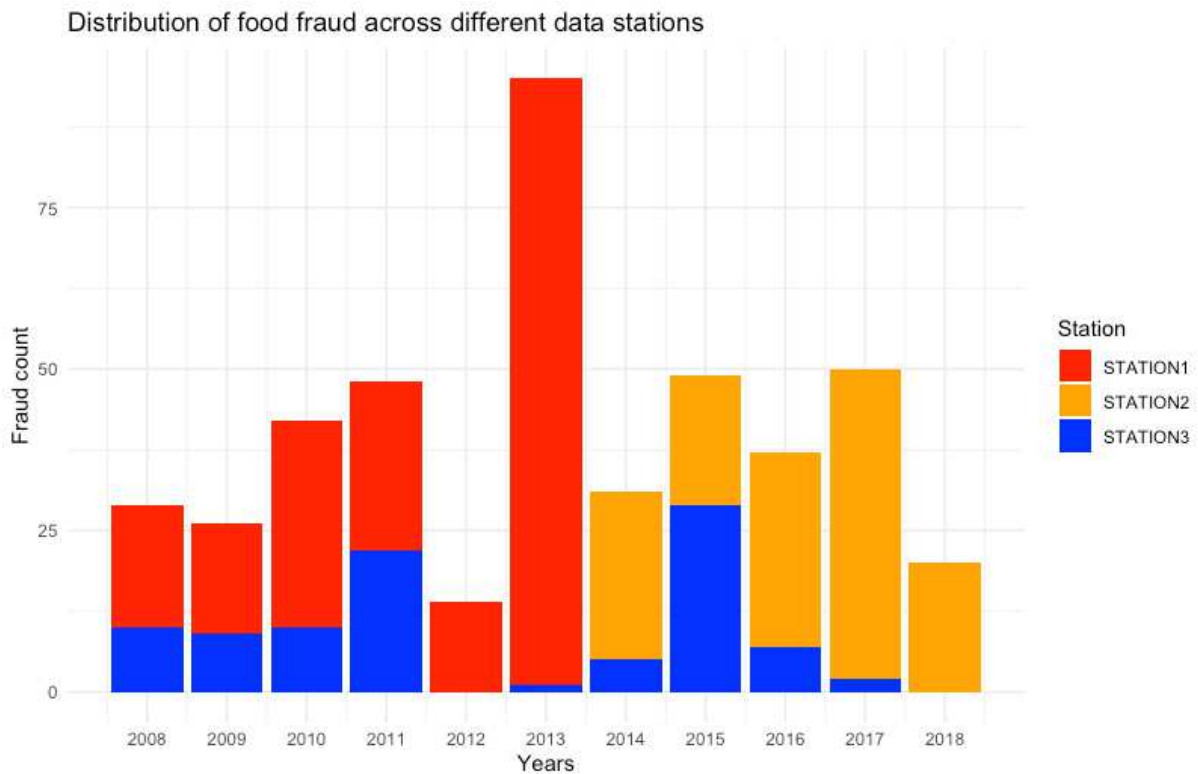
239 **Table 2. Food fraud data available at different data stations**

Data station	Fraud type	Years	Nr of cases	% per data station
STATION-1	Smuggling-Mislabelling-Origin Masking	2008-	105	52%
	Substitution-Dilution	2013	97	48%
STATION-2	Artificial enhancement/Improvement	2014-	1	1%
	Smuggling-Mislabelling-Origin Masking	2018	135	94%

	Substitution-Dilution		8	6%
STATION-3	Artificial enhancement/Improvement	2008- 2018	21	22%
	Smuggling-Mislabelling-Origin Masking		23	24%
	Substitution-Dilution		51	54%

240

241



242

243 **Fig.4.** Summary of data stations depicting imbalance in different fraud types on counts.

244

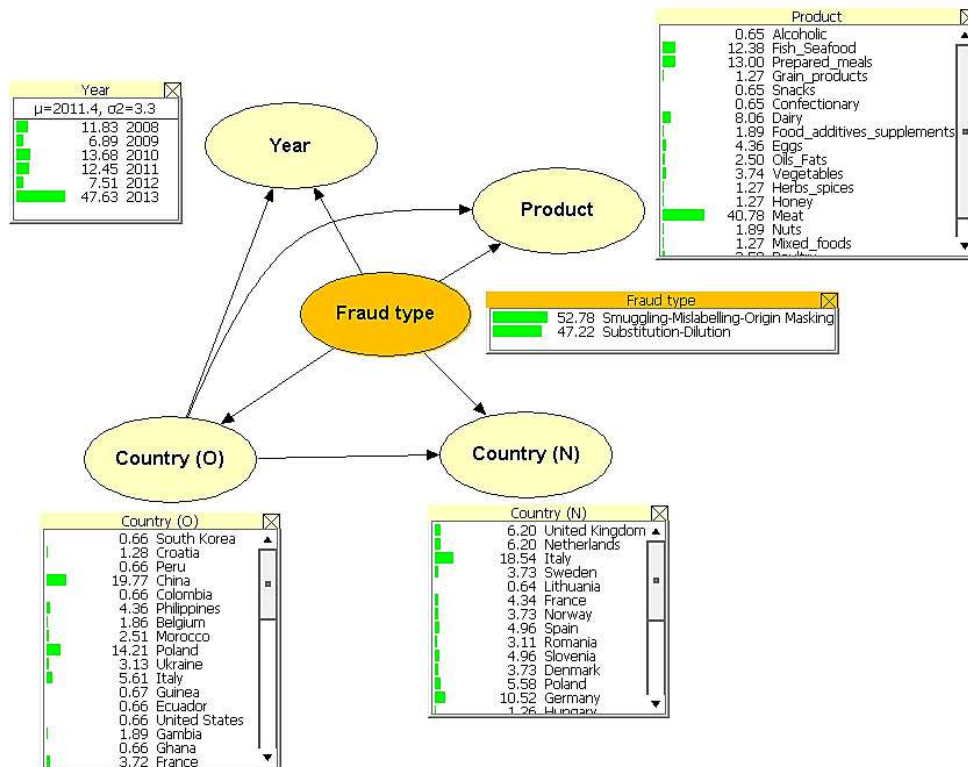
245 **Federated BN**

246 Within the federated learning infrastructure developed, a BN model was trained and validated
 247 to show that knowledge can be shared without source data leaving the data stations and that
 248 such sharing will benefit the decision making of the individual data owner (see experiment 1).
 249 To show that a BN model developed in such federated setting performs like a BN model
 250 developed on pooled data (i.e., traditional manner), experiment 2 was conducted.

251

252 *Experiment 1.*

253 A BN model was created, trained, and tested on each individual data station (i.e., individual BN
254 model) (Table 3) (Fig.5). Its performance was compared to the performance of a combined BN
255 model trained on training data from all data stations (i.e., combined training datasets from all
256 data stations) and then tested on test data from the individual data stations.



257

258 **Fig.5.** Example of the structure of the BN model for STATION-1. The nodes (ellipses in the
259 figure) represent the indicators. The arrows indicate linkages between these nodes. The states
260 are depicted as squares below the nodes.

261

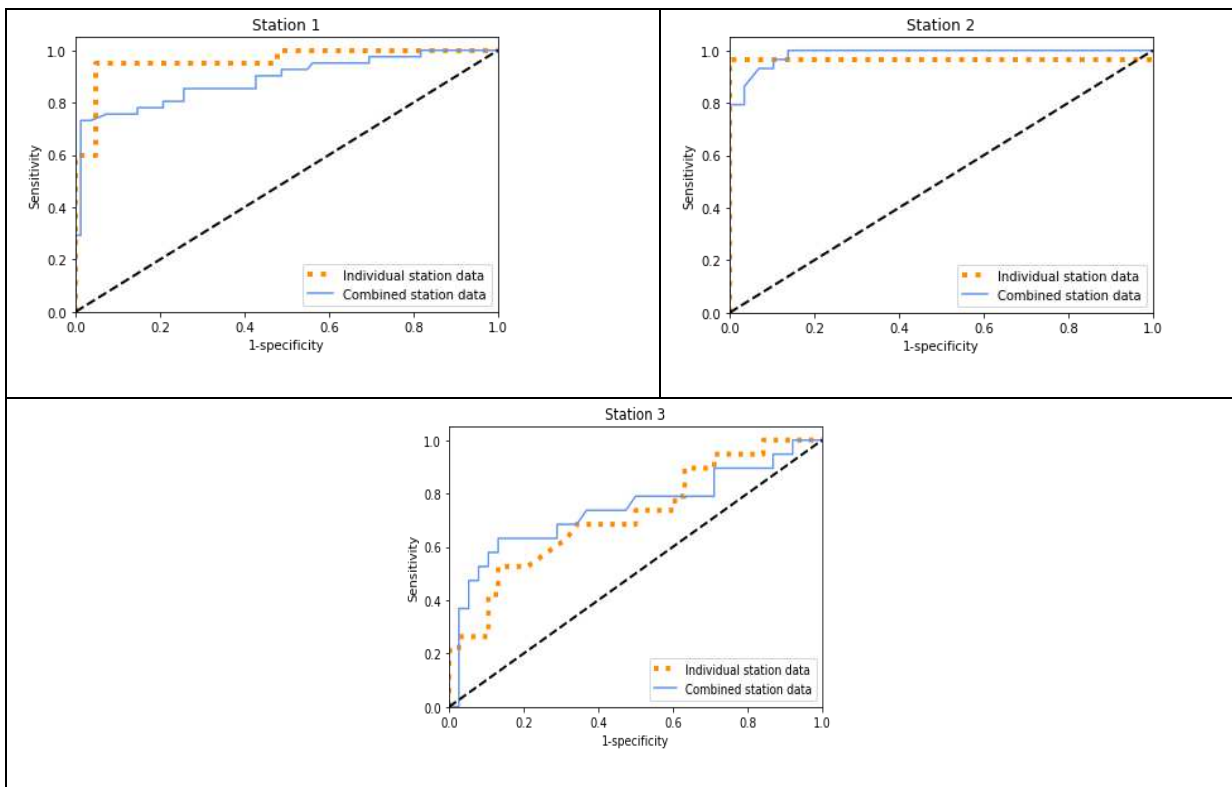
262 ROC curves used to evaluate and compare the performance of the BN models based on the
263 individual station data versus the combined station data for STATION-1, 2, and 3, respectively,
264 are shown in Table 3 and Fig 6. The overall model performance is expressed by the area under
265 the ROC curve (AUC). The results show that the accuracy of each BN model was high for

266 STATION-1 & 2 (i.e., AUC= 0.96), but for STATION-3 the accuracy was significantly lower
267 (i.e., AUC=0.72; Table 3). For the combined BN model when tested with the test data sets from
268 the individual data stations the accuracy was equal to 0.89 for STATION-1, 0.99 for STATION-
269 2 and 0.74 for STATION-3 (Table 3). This shows that a larger number of data available to
270 train/develop a BN model does not necessarily lead to higher accuracies, which can be
271 explained in our case with the data heterogeneity in relation to food fraud type (Fig.4). The
272 relatively small decrease in accuracy of the combined BN compared to the individual BN in
273 STATION-1 (AUC= 0.96 vs. AUC=0.89) is noteworthy because the combined BN model
274 includes three categories of food fraud, whereas the STATION-1 training data set includes only
275 two (see Table 2). However, one should realize that the combined BN model contains more
276 knowledge because it was trained on a broader dataset of food fraud cases than the individual
277 datasets (different products and/or countries of origin) and therefore covers the real situation
278 better and therefore allows the user to make better decisions. In the case presented in this study,
279 the owner of STATION-1 lacked the food category "artificial refinement/improvement" in the
280 dataset, but with the "combined" BN model, the user of STATION-1 effectively gains
281 knowledge about this type of fraud. Moreover, the ROC curves show the trade-off between
282 sensitivity and specificity. An improvement of the combined BN model compared to the
283 individual BN model was observed for the sensitivity parameter of performance, especially for
284 STATION-2 (increase from 0.49 to 0.78, see Table 3). Higher sensitivity means that the
285 combined BN model is better able to identify the food fraud cases. Nevertheless, lower
286 specificity was also found for the combined BN model in STATION-2 (decrease from 0.83 to
287 0.69, see Table 3), which means that the combined BN model leads to more misclassifications
288 of positive food fraud cases where no food fraud is present compared to the single BN model.
289
290

291 **Table 3. AUC, Average sensitivity, Average specificity**

	STATION-1		STATION-2		STATION-3	
	Individual	Combined	Individual	Combined	Individual	Combined
AUC	0.96	0.89	0.96	0.99	0.72	0.74
Average sensitivity	0.75	0.75	0.49	0.78	0.62	0.66
Average specificity	0.69	0.63	0.83	0.69	0.58	0.59

292



293

294 **Fig 6.** Performance of the BN models in each data station (Orange) and in combined station
 295 data (Blue).
 296

297 *Experiment 2.*

298 In this setting, a BN was developed on the total dataset without using a FL infrastructure, hence
 299 the data is shared in a traditional manner (i.e., random split of the dataset, 80% for training and
 300 20% for testing). We conducted this experiment to understand the differences between these

301 two approaches. As shown in Table 4 and Fig.7, an AUC of this BN is 0.86 with an average
302 sensitivity of 0.72 and an average specificity of 0.67.

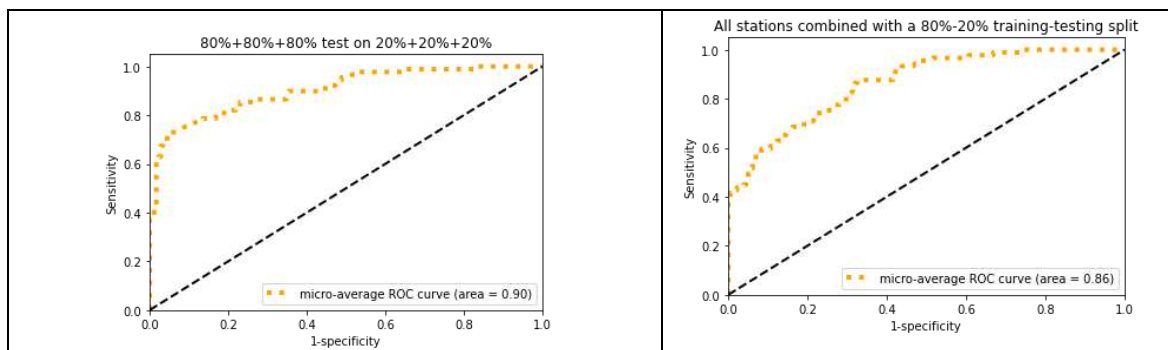
303 As can be seen in Table 4, the combined AUC values with and without a FL infrastructure are
304 very close, with an average AUC value of 0.9 for all three data stations in a federated
305 environment compared to 0.86 without a FL infrastructure. The results show that a FL
306 infrastructure is robust to both model performance and imbalances in the data.

307 **Table 4. AUC, Average sensitivity, Average specificity**

	BN (80%, 20%)	BN (FL datasets combined)
AUC	0.86	0.90
Average sensitivity	0.72	0.77
Average specificity	0.67	0.64

308

309



310

311 **Fig.7.** Performance of the BN models with ROC. Left box: BN Model learned from the
312 combined data in a federated manner. Right box: BN model learned from the combined data
313 in a federated manner.

314

315 An important goal of government agencies, enterprises, and research institutes using a federated
316 infrastructure is to ensure the security and privacy of their data comply with the General Data
317 Protection Regulation (GDPR) and meet data protection compliance measures. In a FL
318 infrastructure, a potential threat is running Docker images on the data stations and granting that
319 permission without the data station owners having visibility into the source code repository of

320 the images. Therefore, it is highly desirable that data station owners allow these images to run
321 on their compute nodes only after validating the source code in the Docker image itself. After
322 validation, the data owner can release this version for execution on their data. The infrastructure
323 handles the validation of the algorithm version using a hash and Docker registry, see the details
324 in (Moncada-Torres, A et al., 2020). Vantage6 uses token-based authentication and
325 authorization for each data station to join the collaboration, while an API can be used to create
326 such collaboration and organizations within that facility.

327 There are currently 3 main issues that are observed in a federated learning setup which are: data
328 democratization, limitation on AI models, and efficiency of the tools. In most federated learning
329 systems, there is not much emphasis on democratization of data, which seems to be of
330 paramount importance for data harmonization. Current constraints in federated learning require
331 that data is in a static structure in a table format (e.g., CSV) and that the order of variables in
332 that data is maintained (Han, et al., 2022; Warnat-Herresthal, et al., 2021). Aside from this, it
333 is important to consider the data type, as it poses a problem when Machine Learning (ML)
334 models use this data (e.g., a ML model might expect a number, but the value is in characters).
335 To address some of these issues, modern data formats have been proposed, such as linked data
336 formats like RDF, where data can be represented as triples. The advantage of this approach is
337 that the data for ML models does not need to be in a predefined tabular structure but can be
338 expanded as it is. However, in order to be able to use these data, an additional SPARQL⁵ layer
339 must be integrated into the machine learning models, which first identifies the required
340 variables and determines whether they are available at all from the data owners. This does not
341 preclude still considering preprocessing steps such as missing values and other data formatting
342 issues. When different data owners collaborate in a federated setting using same data in
343 different contexts additional layers can be added on top of data stations where automated

⁵ <https://www.w3.org/TR/rdf-sparql-query/>

344 deciphering of data based on ontologies can be carried out using modern data standards like
345 FHIR⁶ or OMOP⁷.

346 So far, there are only a limited number of AI models implemented in the FL environment. Some
347 of them are commercial in nature such as HPE⁸, other open-source models that have just been
348 made available in a federated environment are glm (Cellamare, et al., 2022). For large datasets,
349 simple models (e.g., summary statistics) are often sufficient because big data often introduce
350 unique statistical challenges, including scalability and storage bottleneck, noise accumulation,
351 spurious correlation, incidental endogeneity, and measurement errors (Fan, et al., 2014). Apart
352 from that, most complex statistical models are designed to run on single devices in a centralized
353 setting. Modern algorithms like deep learning models need to be redesigned to leverage the
354 power of a FL setup (Chang, et al., 2018).

355 Most ML models are created using languages such as Python or R. These languages allow a
356 researcher to quickly create models for research purposes. However, these models are difficult
357 to operationalize because they have issues with data structures, as most of them work with
358 tabular data and are incompatible with web data formats such as JSON. While there are some
359 packages in these languages that take care of some of the problems, they are not inherently
360 efficient. Most of the models created in a FL environment are dockerized. Docker⁹ provides an
361 environment that allows a ML model to be reusable and reproducible by considering all the
362 dependencies that a ML model requires. However, since both Python and R are interpreted
363 languages, docker images created with these languages are very large, resource intensive, and
364 require good network bandwidth and CPU resources. Recently, efforts are being made to create
365 models in modern compiled languages such as Golang¹⁰ and Rust¹¹. Since the ML models

⁶ <https://www.hl7.org/fhir/overview.html>

⁷ <https://www.ohdsi.org/data-standardization/the-common-data-model/>

⁸ <https://www.hpe.com/us/en/solutions/artificial-intelligence/swarm-learning.html>

⁹ <https://www.docker.com/>

¹⁰ <https://go.dev/>

¹¹ <https://www.rust-lang.org/>

366 created using these tools are in binary format, all dependencies are included in them, making
367 these models more efficient both CPU and in terms of network bandwidth. In future, it is
368 important that more ML models are built using compiled languages.

369

370 **Conclusion**

371 In this study a proof of concept of federated learning approach was demonstrated for the first
372 time for the food safety domain using food fraud as a use case. A new federated BN model was
373 implemented in a federated setting that could be trained on the combined data of databases from
374 geographically different locations, without the source data ever leaving the data stations. It is
375 for the first time that such a principle was demonstrated for three data stations and food fraud,
376 although many more data stations can easily be linked to this infrastructure. The developed
377 approach is applicable to any food safety hazard.

378 The federated learning may help to develop powerful prediction models for the benefit to all
379 actors in the food supply chain while the data will not leave the database of the data owner,
380 hence solving GDPR and business sensitivity issues. Such a concept may stimulate the
381 collaboration along the food supply chain and lead to an increased trust among actors. In
382 addition, making use of data from many stakeholders may also stimulate a more efficient use
383 of resources and reduce the costs of data collection (i.e., food safety monitoring).

384 **Supplement**

385 Source code:

- 386 a. Vantage6 library: <https://vantage6.ai>.
- 387 b. Bayesian Network Algorithm: <https://zenodo.org/record/7394279#.Y4t8nezMKMI>
- 388 c. Docker image: harbor2.vantage6.ai/wur/vtg.wur

389

390

- 392 Almalki, F. A. (2020). Utilizing Drone for Food Quality and Safety Detection using Wireless Sensors. *2020 IEEE*
393 *3rd International Conference on Information Communication and Signal Processing (ICICSP)*, 405–412.
394 <https://doi.org/10.1109/ICICSP50920.2020.9232046>
- 395 Beyan, O., Choudhury, A., van Soest, J., Kohlbacher, O., Zimmermann, L., Stenzhorn, H., Karim, M. R.,
396 Dumontier, M., Decker, S., da Silva Santos, L. O. B., & Dekker, A. (2020). Distributed Analytics on
397 Sensitive Medical Data: The Personal Health Train. *Data Intelligence*, 2(1–2), 96–107.
398 https://doi.org/10.1162/dint_a_00032
- 399 Berkum, v. S., Dengerink, J., Ruben, R. (2018). The food systems approach: sustainable solutions for a sufficient
400 supply of healthy food. Wageningen Economic Research memorandum, 29m 2018-064.
401 <https://doi.org/10.18174/451505>
- 402 Bouzembrak, Y., Klüche, M., Gavai, A., & Marvin, H. J. P. (2019). Internet of Things in food safety: Literature
403 review and a bibliometric analysis. *Trends in Food Science & Technology*, 94, 54–64.
404 <https://doi.org/10.1016/j.tifs.2019.11.002>
- 405 Bouzembrak, Y., & Marvin, H. J. P. (2016). Prediction of food fraud type using data from Rapid Alert System for
406 Food and Feed (RASFF) and Bayesian network modelling. *Food Control*, 61, 180–187.
407 <https://doi.org/10.1016/j.foodcont.2015.09.026>
- 408 Chang, K., Niranjana, B., Carson, L., Yi, D., Brown, J., Beers, A., Rosen, B., R. L.D., Kalpathy-Cramer, J. (2018).
409 Distributed deep learning networks among institutions for medical imaging. *J Am Med Inform Assoc*, 2018
410 Aug; 25(8):945-954. <http://doi.org/10.1093/jamia/ocy017>
- 411 Cellamare, M., Gestel, vVJA., Alradhi, H., Martin, F., & Moncada-Torres, A(2022). A Federated Generalized
412 Linear Model for Privacy-Preserving Analysis. *Algorithms*, 15(7), 243. <https://doi.org/10.3390/a15070243>
- 413 Curry, E. (2016). *The Big Data Value Chain: Definitions, Concepts, and Theoretical Approaches* (W. W. José
414 María Cavanillas, Edward Curry (ed.)). https://link.springer.com/chapter/10.1007/978-3-319-21569-3_3
- 415 Dayan, I., Roth, H. R., Zhong, A., Harouni, A., Gentili, A., Abidin, A. Z., Liu, A., Costa, A. B., Wood, B. J., Tsai,
416 C.-S., Wang, C.-H., Hsu, C.-N., Lee, C. K., Ruan, P., Xu, D., Wu, D., Huang, E., Kitamura, F. C., Lacey,
417 G., ... Li, Q. (2021). Federated learning for predicting clinical outcomes in patients with COVID-19. *Nature*
418 *Medicine*, 27(10), 1735–1743. <https://doi.org/10.1038/s41591-021-01506-3>
- 419 Deist, T. M., Dankers, F. J. W. M., Ojha, P., Scott Marshall, M., Janssen, T., Faivre-Finn, C., Masciocchi, C.,
420 Valentini, V., Wang, J., Chen, J., Zhang, Z., Spezi, E., Button, M., Jan Nuytens, J., Vernhout, R., van Soest,
421 J., Jochems, A., Monshouwer, R., Bussink, J., ... Dekker, A. (2020). Distributed learning on 20 000+ lung
422 cancer patients – The Personal Health Train. *Radiotherapy and Oncology*, 144, 189–200.
423 <https://doi.org/10.1016/j.radonc.2019.11.019>
- 424 Dooley, D. M., Griffiths, E. J., Gosal, G. S., Buttigieg, P. L., Hoehndorf, R., Lange, M. C., Schriml, L. M.,
425 Brinkman, F. S. L., & Hsiao, W. W. L. (2018). FoodOn: a harmonized food ontology to increase global food
426 traceability, quality control and data integration. *Npj Science of Food*, 2(1), 23.
427 <https://doi.org/10.1038/s41538-018-0032-6>
- 428 Fan, Jianqing., Han, Fang., & Liu, Han(2014). Challenges of Big Data Analysis. *Natl Sci Rev*. 2014 June; 1(2):
429 293-314. <https://doi.org/10.1093/nsr/nwt032>
- 430 Flores, M., Dayan, I., Roth, H., Zhong, A., Harouni, A., Gentili, A., Abidin, A., Liu, A., Costa, A., Wood, B., Tsai,
431 C.-S., Wang, C.-H., Hsu, C.-N., Lee, C. K., Ruan, C., Xu, D., Wu, D., Huang, E., Kitamura, F., ... Wen, Y.
432 (2021). Federated Learning used for predicting outcomes in SARS-COV-2 patients. *Research Square*.
433 <https://doi.org/10.21203/rs.3.rs-126892/v1>
- 434 Geleijnse, G., Chiang, R. C.-J., Sieswerda, M., Schuurman, M., Lee, K. C., van Soest, J., Dekker, A., Lee, W.-C.,
435 & Verbeek, X. A. A. M. (2020). Prognostic factors analysis for oral cavity cancer survival in the Netherlands
436 and Taiwan using a privacy-preserving federated infrastructure. *Scientific Reports*, 10(1), 20526.
437 <https://doi.org/10.1038/s41598-020-77476-2>
- 438 Heringa, J., Dumon, O., van der Lei, J., Sansone, S.-A., Brookes, A. J., Dillo, I., da Silva Santos, L. B., Grethe, J.
439 S., Goble, C., Waagmeester, A., Dumontier, M., Roos, M., Slater, T., Clark, T., Zhao, J., Bouwman, J.,
440 Rocca-Serra, P., Swertz, M. A., Wittenburg, P., ... Wilkinson, M. D. (2016). The FAIR Guiding Principles
441 for scientific data management and stewardship. *Scientific Data*. <https://doi.org/10.1038/sdata.2016.18>
- 442 Jin, C., Bouzembrak, Y., Zhou, J., Liang, Q., van den Bulk, L. M., Gavai, A., Liu, N., van den Heuvel, L. J.,
443 Hoenderdaal, W., & Marvin, H. J. P. (2020). Big Data in food safety- A review. *Current Opinion in Food*
444 *Science*, 36, 24–32. <https://doi.org/10.1016/j.cofs.2020.11.006>
- 445 Johan van Soest, C. S. O. M. M. P. B. van den B. A. M. C. van O. D. T. A. D. M. D. (2018). Using the Personal
446 Health Train for Automated and Privacy-Preserving Analytics on Vertically Partitioned Data. In *Building*
447 *Continents of Knowledge in Oceans of Data: The Future of Co-Created eHealth* (Vol. 247, pp. 581–585).
448
- 449 Lokhorst, C., de Mol, R. M., & Kamphuis, C. (2019). Invited review: Big Data in precision dairy farming. *Animal*,

450 13(7), 1519–1528. <https://doi.org/10.1017/S1751731118003439>

451 Moncada-Torres, A., & Frank Martin I, Melle Sieswerda, Johan Van Soest, G. G. (n.d.). *VANTAGE6: an open*

452 *source priVAcY preserviNg federaTed leArninG infrastructurE for Secure Insight eXchange.*

453 <https://pubmed.ncbi.nlm.nih.gov/33936462>

454 Nelis, J. L. D., Tsagkaris, A. S., Dillon, M. J., Hajslova, J., & Elliott, C. T. (2020). Smartphone-based optical

455 assays in the food safety field. *TrAC Trends in Analytical Chemistry*, 129, 115934.

456 <https://doi.org/10.1016/j.trac.2020.115934>

457 Scutari, M. (2010). Learning Bayesian Networks with the bnlearn R Package. *Journal of Statistical Software*, 2.

458 Shi, Z., Zhovannik, I., Traverso, A., Dankers, F. J. W. M., Deist, T. M., Kalendralis, P., Monshouwer, R., Bussink,

459 J., Fijten, R., Aerts, H. J. W. L., Dekker, A., & Wee, L. (2019). Distributed radiomics as a signature

460 validation study using the Personal Health Train infrastructure. *Scientific Data*, 6(1), 218.

461 <https://doi.org/10.1038/s41597-019-0241-0>

462 Marvin, Hans J.P., Yamine Bouzembrak, H.J. van der Fels-Klerx, Corné Kempenaar, Roel Veerkamp, Aneesh

463 Chauhan, Sanne Stroosnijder, Jan Top, Gökem Simsek-Senel, Hans Vrolijk, Willem Jan Knibbe, Lu Zhang,

464 Remko Boom, Bedir Tekinerdogan. (2022). Digitalisation and Artificial Intelligence for sustainable food

465 systems. *Trends in Food Science & Technology* 120 (2022) 344–348.

466 <https://doi.org/10.1016/j.tifs.2022.01.020>

467 Hans J.P. Marvin, Esmée M. Janssen, Yamine Bouzembrak, Peter J.M. Hendriksen and Martijn Staats. (2017).

468 Big data in food safety; an overview. *Critical Reviews in Food Science and Nutrition*, 57:11, 2286-2295,

469 DOI: 10.1080/10408398.2016.1257481

470 H.J.P. Marvin, Y. Bouzembrak, E.M. Janssen, H.J. van der Fels Klerx, E.D. van Asselt & G.A. Kleter (2016). A

471 holistic approach to food safety risks: Food fraud as an example. *Food Research International* 89: 463–470. DOI:

472 10.1016/j.foodres.2016.08.028

473 Hans J. P. Marvin & Y. Bouzembrak (2020). A system approach towards prediction of food safety hazards: Impact

474 of climate and agrichemical use on the occurrence of food safety hazards. *Agricultural Systems*, Volume 178,

475 February 2020, 102760. <https://doi.org/10.1016/j.agsy.2019.102760>

476 Top, J. S. Janssen, H. Boogaard, R. Knapen, Gorkem Simsek-Senel (2022) Cultivating FAIR principles for agri-

477 food data. *Computers and Electronics in Agriculture* 196 (2022) 106909

478 Xiyao Wang, Yamine Bouzembrak, Hans J.P. Marvin, Dave Clarke, Francis Butler (2022) Bayesian Networks

479 modelling of Diarrhetic Shellfish Poisoning in *Mytilus edulis* harvested in Bantry Bay, Ireland. *Harmful*

480 *Algae* 112 (2022) 102171. <https://doi.org/10.1016/j.hal.2021.102171>

481 da Silva Santos, L. O. B., Burger, K., Kaliyaperumal, R., & Wilkinson, M. D. (2022). FAIR Data Point: A

482 FAIR-Oriented Approach for Metadata Publication. *Data Intelligence*, 1-21.

483 Han, J., Ma, Y. F., Han, Y., Zhang, Y., & Huang, G. (2022). Demystifying Swarm Learning: A New Paradigm

484 of Blockchain-based Decentralized Federated Learning. In

485 Hard, A., Rao, K., Mathews, R., Ramaswamy, S., Beaufays, F., Augenstein, S., Eichner, H., Kiddon, C., &

486 Ramage, D. (2018). Federated learning for mobile keyboard prediction. arXiv preprint

487 arXiv:1811.03604.

488 Li, C., Niu, D., Jiang, B., Zuo, X., & Yang, J. (2021). Meta-har: Federated representation learning for human

489 activity recognition. In *Proceedings of the Web Conference 2021* (pp. 912-922).

490 Ma, J., Zhang, Q., Lou, J., Xiong, L., & Ho, J. C. (2021). Communication efficient federated generalized tensor

491 factorization for collaborative health data analytics. In *Proceedings of the Web Conference 2021* (pp.

492 171-182).

493 Warnat-Herresthal, S., Schultze, H., Shastry, K. L., Manamohan, S., Mukherjee, S., Garg, V., Sarveswara, R.,

494 Händler, K., Pickkers, P., Aziz, N. A., Ktena, S., Tran, F., Bitzer, M., Ossowski, S., Casadei, N., Herr,

495 C., Petersheim, D., Behrends, U., Kern, F., Fehlmann, T., Schommers, P., Lehmann, C., Augustin, M.,

496 Rybniker, J., Altmüller, J., Mishra, N., Bernardes, J. P., Krämer, B., Bonaguro, L., Schulte-Schrepping,

497 J., De Domenico, E., Siever, C., Kraut, M., Desai, M., Monnet, B., Saridaki, M., Siegel, C. M., Drews,

498 A., Nuesch-Germano, M., Theis, H., Heyckendorf, J., Schreiber, S., Kim-Hellmuth, S., Balfanz, P.,

499 Eggermann, T., Boor, P., Hausmann, R., Kuhn, H., Isfort, S., Stingl, J. C., Schmalzing, G., Kuhl, C. K.,

500 Röhrig, R., Marx, G., Uhlig, S., Dahl, E., Müller-Wieland, D., Dreher, M., Marx, N., Nattermann, J.,

501 Skowasch, D., Kurth, I., Keller, A., Bals, R., Nürnberg, P., Rieß, O., Rosenstiel, P., Netea, M. G., Theis,

502 F., Mukherjee, S., Backes, M., Aschenbrenner, A. C., Ulas, T., Angelov, A., Bartholomäus, A., Becker,

503 A., Bezdan, D., Blumert, C., Bonifacio, E., Bork, P., Boyke, B., Blum, H., Clavel, T., Colome-Tatche,

504 M., Cornberg, M., De La Rosa Velázquez, I. A., Diefenbach, A., Diltney, A., Fischer, N., Förstner, K.,

505 Franzenburg, S., Frick, J.-S., Gabernet, G., Gagneur, J., Ganzenmueller, T., Gauder, M., Geißert, J.,

506 Goesmann, A., Göpel, S., Grundhoff, A., Grundmann, H., Hain, T., Hanses, F., Hehr, U., Heimbach, A.,

507 Hoeper, M., Horn, F., Hübschmann, D., Hummel, M., Iftner, T., Iftner, A., Illig, T., Janssen, S.,

508 Kalinowski, J., Kallies, R., Kehr, B., Keppler, O. T., Klein, C., Knop, M., Kohlbacher, O., Köhrer, K.,

509 Korbel, J., Kremsner, P. G., Kühnert, D., Landthaler, M., Li, Y., Ludwig, K. U., Makarewicz, O., Marz,

510 M., McHardy, A. C., Mertes, C., Münchhoff, M., Nahsen, S., Nöthen, M., Ntoumi, F., Overmann, J.,
511 Peter, S., Pfeffer, K., Pink, I., Poetsch, A. R., Protzer, U., Pühler, A., Rajewsky, N., Ralser, M., Reiche,
512 K., Ripke, S., da Rocha, U. N., Saliba, A.-E., Sander, L. E., Sawitzki, B., Scheithauer, S., Schiffer, P.,
513 Schmid-Burgk, J., Schneider, W., Schulte, E.-C., Sczyrba, A., Sharaf, M. L., Singh, Y., Sonnabend, M.,
514 Stegle, O., Stoye, J., Vehreschild, J., Velavan, T. P., Vogel, J., Volland, S., von Kleist, M., Walker, A.,
515 Walter, J., Wieczorek, D., Winkler, S., Ziebuhr, J., Breteler, M. M. B., Giamarellos-Bourboulis, E. J.,
516 Kox, M., Becker, M., Cheran, S., Woodacre, M. S., Goh, E. L., Schultze, J. L., Study, C.-A., &
517 Deutsche, C.-O. I. (2021). Swarm Learning for decentralized and confidential clinical machine learning.
518 *Nature*, 594, 265-270.

519 Wu, J., Liu, Q., Huang, Z., Ning, Y., Wang, H., Chen, E., Yi, J., & Zhou, B. (2021). Hierarchical personalized
520 federated learning for user modeling. In *Proceedings of the Web Conference 2021* (pp. 957-968).
521