

## Journal Pre-proofs

Intestinal microbiota alterations by dietary exposure to chemicals from food cooking and processing. Application of Data Science for risk prediction

Sergio Ruiz-Saavedra, Herminio García-González, Silvia Arboleya, Nuria Salazar, José Emilio Labra-Gayo, Irene Díaz, Miguel Gueimonde, Sonia González, Clara G. de los Reyes-Gavilán

PII: S2001-0370(21)00041-6  
DOI: <https://doi.org/10.1016/j.csbj.2021.01.037>  
Reference: CSBJ 857

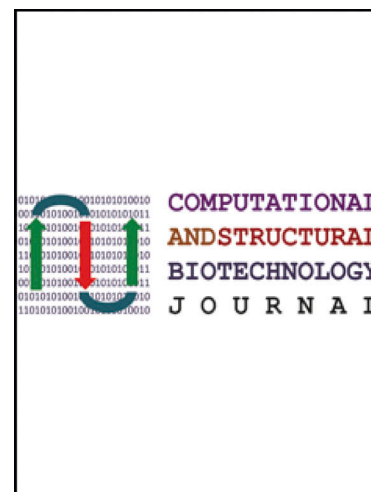
To appear in: *Computational and Structural Biotechnology Journal*

Received Date: 28 August 2020  
Revised Date: 22 January 2021  
Accepted Date: 22 January 2021

Please cite this article as: S. Ruiz-Saavedra, H. García-González, S. Arboleya, N. Salazar, J. Emilio Labra-Gayo, I. Díaz, M. Gueimonde, S. González, C.G. de los Reyes-Gavilán, Intestinal microbiota alterations by dietary exposure to chemicals from food cooking and processing. Application of Data Science for risk prediction, *Computational and Structural Biotechnology Journal* (2021), doi: <https://doi.org/10.1016/j.csbj.2021.01.037>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2021 The Author(s). Published by Elsevier B.V. on behalf of Research Network of Computational and Structural Biotechnology.



1 Title: **Intestinal microbiota alterations by dietary exposure to chemicals from food**  
2 **cooking and processing. Application of Data Science for risk prediction**

3  
4  
5 Authors: Sergio Ruiz-Saavedra<sup>a,b,c,†</sup>, Herminio García-González<sup>d,e,†</sup>, Silvia Arboleya<sup>a,c</sup>,  
6 Nuria Salazar<sup>a,c</sup>, José Emilio Labra-Gayo<sup>d</sup>, Irene Díaz<sup>d</sup>, Miguel Gueimonde<sup>a,c</sup>, Sonia  
7 González<sup>b,c</sup> and Clara G. de los Reyes-Gavilán<sup>a,c\*</sup>

8  
9 Affiliations

- 10  
11 a. *Department of Microbiology and Biochemistry of Dairy Products, Instituto de*  
12 *Productos Lácteos de Asturias (IPLA-CSIC), 33300 Villaviciosa, Asturias, Spain.*  
13 b. *Department of Functional Biology, University of Oviedo, 33006 Oviedo, Asturias.*  
14 *Spain.*  
15 c. *Diet, Microbiota and Health Group, Instituto de Investigación Sanitaria del*  
16 *Principado de Asturias (ISPA), 33011 Oviedo, Spain.*  
17 d. *Department of Computer Science, University of Oviedo, C/ Federico García Lorca*  
18 *S/N 33007 Oviedo, Asturias, Spain.*  
19 e. *IT and Communications Service, University of Oviedo, C/ Fernando Bongera S/N*  
20 *33006 Oviedo, Asturias, Spain.*  
21 †. *Both authors contributed equally to this work.*

22  
23 \* Corresponding author

24  
25 Email address: [greyes\\_gavilan@ipla.csic.es](mailto:greyes_gavilan@ipla.csic.es)

26 Phone: +34985893335

27  
28  
29  
30  
31  
32  
33  
34  
35  
36 *Abbreviations:* ANN, Artificial Neural Networks; BaP, Benzo(a)pyrene; CHARRED,  
37 Computerized Heterocyclic Amines Resource for Research Epidemiology of Disease;  
38 CRC, Colorectal Cancer; DT, Decision Tree; EPIC, European Prospective Investigation  
39 into Cancer and Nutrition; FFQ, Food Frequency Questionnaire; HCA, Heterocyclic  
40 Aromatic Amines; IARC, International Agency for Research on Cancer; IM, Intestinal  
41 Microbiota; KNN, k-Nearest Neighbour; miRNAs, micro-RNAs; NA, Nitrosamines;  
42 NIH-AARP, National Institute of Health-American Association of Retired Persons;  
43 PAH, Polycyclic Aromatic Hydrocarbons; RF, Random Forest; RDF, Resource  
44 Description Framework; SPARQL, Protocol and RDF Query Language; SVM, Support  
45 Vector Machine; WHO, World Health Organization.

46  
47  
48  
49 **Abstract**

50 Diet is one of the main sources of exposure to **toxic chemicals with carcinogenic**  
51 **potential**, **some of which** are generated during food processing, depending on the type of  
52 food (**primarily** meat, fish, **bread and potatoes**), cooking methods and temperature.  
53 Although demonstrated in animal models at high doses, an unequivocal link between  
54 dietary **exposure to these compounds** with disease has not been proven in humans. A  
55 major difficulty in assessing the actual intake of **these toxic compounds** is the lack of  
56 standardised and harmonised protocols for collecting and analysing dietary information.  
57 The intestinal microbiota (IM) has a great influence on health and is altered in some  
58 diseases **such** as colorectal cancer (CRC). Diet influences **the** composition and activity  
59 of the IM, and the net exposure to genotoxicity of **potential** dietary carcinogens in the  
60 gut depends on the interaction among these compounds, IM and diet. This review  
61 analyses critically **the** difficulties and challenges in the study of interactions among  
62 these three actors on the onset of CRC. Machine Learning (ML) of data obtained in  
63 subclinical and precancerous stages would help to establish risk thresholds for the  
64 intake of **toxic compounds generated during food processing** as related to diet and IM  
65 profiles, whereas Semantic Web could improve data accessibility and usability from  
66 different studies, as well as **helping to elucidate** novel interactions among **those**  
67 **chemicals**, IM and diet.

68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98

Keywords: intestinal microbiota, colorectal cancer, diet, **toxic chemicals**, machine learning, semantic web

## 99 1. Introduction

100 Diet is one of the main **sources** of exposure to **toxic** compounds **with carcinogenic**  
101 **potential**. In October 2015 the International Agency for Research on Cancer **from the**  
102 **World Health Organization** (IARC-WHO) announced the classification of processed  
103 meat as “carcinogenic to humans” and red meat as “probably carcinogenic to humans”  
104 [1]. Diets from most developed countries are characterized by high intakes of meat,  
105 **which** is often fried, griddled or barbecued, and by an increasing consumption of  
106 processed foods. When cooking muscle meat from animals or fish at high temperature,  
107 some **chemicals** are formed at levels **that** depend on the cooking procedure and  
108 temperature; **some of** these compounds **can** cause cancer when administered at high  
109 doses in experimental animals [2]. However, although the intake of dietary **compounds**  
110 **with carcinogenic potential** in humans is considerably lower than in experimental  
111 animals, **lifetime** exposure can differ considerably among individuals. No regulations  
112 exist about the presence in foods of cooking -related **potential carcinogens**. This **aspect**  
113 is specially relevant for public health, as most cooking mutagen/**genotoxic compounds**  
114 are generated at home, restaurants and local ready-to-eat food providers.

115 **Despite** that some international projects have evaluated the association between  
116 nutrition (including cooking methods) and cancer, such as the *European Prospective*  
117 *Investigation into Cancer and Nutrition* (EPIC) or the NIH-AARP Diet and Health  
118 Study, an unequivocal link between dietary exposure to **chemicals** and human cancer [3]  
119 has not been **shown**. **The** underlying reasons for this may be **as follows**: i) the difficulty  
120 to determine the exact exposure to **these compounds** (depending not only on the intake  
121 but **also** on the cumulative exposure and delayed effect through life), ii)  
122 interindividual variation in the detoxifying activity of endogenous enzymes, iii)  
123 cumulative exposure to toxic **compounds** from different environmental sources, iv)  
124 synergistic interaction among different **compounds** and, v) the role, not sufficiently  
125 explored **to date**, of the interaction between diet and the intestinal microbiota (IM) on  
126 the net carcinogenic potential. Therefore, studies designed to explore these interactions  
127 could help to establish risk thresholds for disease as a function of dietary intake **of**  
128 **potential carcinogens**, global diet and microbiota. The present review analyses  
129 difficulties inherent to this type of studies and how Machine Learning (ML) and  
130 Semantic Web could assist in data modelling for risk assessment.

131

## 132 2. **Chemicals with carcinogenic potential** formed during food cooking and 133 **processing**

134 One of the most important risk factors for the development of cancer is the exposure to  
135 dietary **toxic chemicals with** carcinogenic **and pro-carcinogenic potential** **which**, when  
136 consumed regularly at certain levels, can increase the risk of triggering tumorigenic  
137 processes. Nitrates, nitrites, nitrosamines (NA), heterocyclic amines (HCA), polycyclic  
138 aromatic hydrocarbons (PAH) **and acrylamide**, are amongst the substances with the  
139 highest carcinogenic potential. **Some of** these compounds are not naturally present in  
140 foods but can be incorporated (nitrates and nitrites) or generated (NA, HCA and PAH)  
141 during the processing of foodstuffs containing nitrogenous and creatine components by  
142 heat-direct exposure procedures [4]. HCA have accumulated solid scientific evidence as  
143 cancer risk factors and are the only carcinogens formed exclusively during the cooking  
144 process. Specifically, HCA show a mutagenicity index more than 1000 times higher  
145 than benzo(a)pyrene (BaP) [3]. Carcinogens may act through various mechanisms, such  
146 as chromosomal aberrations, single strand breaks and DNA adducts or **oestrogenic**

147 activity [5]. Several prospective cohort studies reported mean intakes of HCA between  
148 69.4 ng/day and 821 ng/day in European countries [6, 7] and from 49.95 ng/day to  
149 151.9 ng/day in Chinese communities and the United States [8, 9]. The observed  
150 variability among countries and individuals may be attributed to differences in the  
151 methodology used for the assessment of potential carcinogenic chemicals and to  
152 differences in dietary patterns and cooking preferences around the world. For example,  
153 compared to the 134.5 ng/day contribution of 50 g of broiled beef (0.00269 ppm/day),  
154 one daily serving of 50 g of broiled chicken could increase the intake of HCAs  
155 (PhIP+MeIQx) by 1350 ng/day (0.027 ppm/day) [10]. Induction of tumours in the large  
156 intestine of F344 rats and C57BL/6 mice have been demonstrated during prolonged  
157 exposure (40 to 72 weeks) to high concentrations of some HCA in diet (i.e. 300  
158 ppm/day) [2]. Although useful to demonstrate tumorigenic potential, experiments with  
159 animals are not intended to predict true human cancer incidence associated with  
160 exposure to chemicals.

161 PAHs are found in cured and processed meat and fats, primarily [11]. Dietary exposure  
162 levels ranged from the order of ng/day in some Asian publications [12] to the order of  
163 µg/day reported in other publications [13]. BaP is the most-used marker to detect the  
164 presence of PAHs in foods [14, 15]. NA are detected in cured meat and smoked foods  
165 and are also endogenously formed from the interaction of nitrosating agents with  
166 amines and amides [16]. The intake of NA showed unclear relationships with  
167 pancreatic-cancer but positive associations with colorectal cancer (CRC) and gastric  
168 cancer [17, 18].

169 Nitrates and nitrites are often used as food additives in processed meats, fish, cheese,  
170 and fermented products, to preserve them from microbial alteration [19]. The  
171 simultaneous presence in certain foods of amino acids can lead to a chemical reaction  
172 that results in the formation of NA, especially when a heat treatment is applied; N-  
173 nitrosopyrrolidine (NPYR) and N-nitrosodimethylamine (NMDA) are the NA most  
174 frequently found in foods [19]. Several studies have shown an increased risk of CRC  
175 development for NMDA intakes of 0.03 - 0.07 µg/day [20].

176 Acrylamide is formed by asparagine decarboxylation in the presence of reducing sugars  
177 during nonenzymatic browning (Maillard reaction) [21]. It is naturally found in foods,  
178 but can also form during the thermal treatment. In European countries, the major  
179 sources of acrylamide are potatoes, coffee and cereal products [22]. Acrylamide has  
180 been classified by the EFSA [23] as probably carcinogenic to humans. However, there  
181 is still no regulation on the maximum recommended intake albeit there is a general  
182 recommendation to limit its consumption.

183

### 184 3. Challenges to determine the actual intake of toxic chemicals with carcinogenic 185 potential generated during food cooking and processing

186 Recent meta-analyses of epidemiological studies are still not completely conclusive  
187 about the relationship of the intake of toxic compounds with carcinogenic potential  
188 resulting from food processing and cancer development [3] as it is complex to  
189 disentangle the effect of these compounds from the effect of the food itself. Most of the  
190 research revealing the impact of red and processed meat consumption in the relative risk  
191 of developing several chronic pathologies, such as CRC, prostate or lung cancer is the  
192 result of longitudinal epidemiological studies. Although these studies are useful from a  
193 descriptive point of view and for the generation of research hypotheses, they have a

194 limited potential for the establishment of cause-effect relationships, leading to the  
195 continuing debate about the health impact of meat intake.

196 A major difficulty in assessing quantitatively the actual intake of food potential  
197 carcinogens in the population is the selection of the most appropriate method for the  
198 collection of dietary data. The food frequency questionnaire (FFQ), multiple day food  
199 records and 24-hour dietary recall are among the most extensively used tools for this  
200 purpose. With independence to the systematic and random errors inherent to these  
201 methods [24], some factors such as the time period covered by the dietary  
202 questionnaires and the number of items included or the quantification of the portions  
203 consumed, affect the quality of the information collected and therefore the conclusions  
204 drawn. It is important to note that the risk of developing cancer from exposure to  
205 environmental factors, including diet and lifestyle, is cumulative over a subject's  
206 lifetime. For this reason, it seems more appropriate to use questionnaires with the  
207 capacity to describe long-term dietary habits, such as the FFQ. However, the FFQ has  
208 the disadvantage of providing less accurate information on energy and nutrient intake  
209 compared with the other methods mentioned above. In addition, some of the postulated  
210 mechanisms linking meat consumption to cancer risk include the content of these foods  
211 in HCA [4], PAH and other compounds generated during the high-temperature  
212 processing of foods, particularly in meats cooked at "well-done" degree [4]. Therefore,  
213 at the time of quantifying the intake of different toxic compounds with carcinogenic  
214 potential, it is important to detail in a harmonised way some characteristics related to  
215 the culinary preparation of foods, such as cooking time, processing method, temperature  
216 or degree of browning [11]. This is a strong add-on difficulty because it prolongs the  
217 duration of the baseline questionnaires, increasing the number of items included. In  
218 addition, the analysis of the information obtained is more complex than usual for the  
219 calculation of a nutrient, since for each of the foods surveyed, the type of processing  
220 (preservation or cooking) and the duration and temperature of cooking should be  
221 considered. The estimation of dietary compounds with carcinogenic potential can be  
222 extracted from information compiled in various databases. The most widely used  
223 databases are those developed by the EPIC study for the European population [25] and  
224 by the Computerized Heterocyclic Amines Resource for Research in Epidemiology of  
225 Disease (CHARRED) database for the United States [26]. Both databases provide key  
226 information for integrating the analysis of dietary potential carcinogens on a systematic  
227 basis. The EPIC database compiles information obtained from 139 references regarding  
228 the content per 100 g of food in NA, HAC, PAH, nitrites and nitrates in more than 200  
229 food items. The food composition table is classified according to the preservation  
230 method, cooking method, degree of browning and temperature [25]. This information is  
231 also present in the CHARRED database, which has developed a special module within a  
232 FFQ in conjunction with the mutagens database to estimate intake of the mutagenic  
233 compounds in cooked meats [26]. In addition, acrylamide content was estimated from  
234 the EFSA categorisation of European food products for monitoring purposes [27].

235 A broader approach is necessary in the future in order to lay the foundations for  
236 improving the understanding of the complex diet-cancer association in the long term.  
237 This approach would require consensus on standardised and harmonised protocols for  
238 collecting dietary information, classifying the degree of cooking and calculating  
239 carcinogens derived from food processing. This method should be complemented with  
240 advanced tools for mathematical analysis of data that enable researchers to both identify  
241 risk factors for these pathologies and explain their impact in the complex context of a  
242 subject's global diet and lifestyles.

243

#### 244 4. Intestinal microbiota and human health. Methods to study composition and 245 functionality

246 The IM is defined as the set of microorganisms inhabiting the intestine. The microbiota  
247 has co-evolved with the host over thousands of years, leading to the establishment of a  
248 mutually beneficial microbiota-host relationship. The number of microorganisms in the  
249 human gut exceeds  $10^{14}$  and this microbiota encodes a collection of genes ~10 times  
250 greater than these encoded by the human genome, providing exclusive capabilities and  
251 functions essential for the maintenance of health. The role of the IM begins in early life,  
252 participating in the development of the host's immune, digestive and nervous systems  
253 by strengthening intestinal epithelium integrity and gut barrier, protecting against  
254 pathogens and playing a major role in helping to harvest nutrients and energy from our  
255 diet. Therefore, the IM results in a key player for host physiology [28].

256 This IM represents a large factory producing bioactive compounds and participating in  
257 the host's metabolism and nutrition. Actually, host metabolism is the combination of the  
258 capabilities of both the human and the IM genomes. The microbiota ferments  
259 indigestible complex carbohydrates and proteins from the diet producing short-chain  
260 fatty acids, primarily acetate, propionate and butyrate, which are quickly absorbed by  
261 the gut epithelial cells [29]. Acetate is primarily delivered to peripheral tissues for use  
262 as a substrate in the synthesis of cholesterol and fatty acids; propionate is absorbed in  
263 the liver and participates in gluconeogenesis; and butyrate is used as one of the main  
264 energy sources by colonocytes. Other metabolites are also produced by the IM such as  
265 branched chain fatty acids, secondary bile acids, amino acids, trimethylamine,  
266 neurotransmitters, and some essential vitamins [30, 31]. Some of these metabolites may  
267 suffer further transformations, such as the case of trimethylamine which, upon  
268 absorption will be oxidised in the liver to trimethylamine-N-Oxide, a known risk factor  
269 for cardiovascular disease. Therefore, all these metabolites participate in the host's  
270 physiology and strong evidence now supports the role of the IM in the maintenance of  
271 human homeostasis. For this reason, adverse changes in the gut microbiota composition  
272 and/or function, the so-called *dysbiosis*, are related to different gastrointestinal  
273 disorders, such as diarrhoea, inflammatory bowel disease, cancer, or extra-intestinal  
274 diseases such as obesity, allergies, neurological sicknesses or other metabolic diseases.  
275 Different stressors, including dietary changes, antibiotic or other drugs treatments, and  
276 carcinogens from the diet can be involved in the development of dysbiosis.

277 Members of Bacteroidetes and Firmicutes phyla followed by Actinobacteria,  
278 Proteobacteria and Verrucomicrobia primarily make up the composition of the adult IM.  
279 However, at lower taxonomical levels, the complexity of the IM is higher and is  
280 represented by thousands of different microbial species. This diversity also occurs  
281 among individuals, making almost impossible the definition of a *normal* or *healthy* IM  
282 composition for an entire population. However, it is also known that the IM exhibits  
283 high functional redundancy, meaning that some functions may be conferred by multiple  
284 bacteria, from related and unrelated species, making the IM more conserved at the  
285 functional than at compositional level [32]. Accounting for this variability, some  
286 authors have tried to define the "normal or healthy" IM as the "intestinal microbial  
287 community that assist the host to maintain a healthy status under certain environmental  
288 conditions" [33], understanding that under different environmental conditions including  
289 dietary habits the optimal microbiota for health may also be different. For this reason,  
290 when we aim to assess the effect of a specific factor or a specific disease on the gut

291 microbiota, it is crucial to identify the specific alterations present in the gut microbiota  
292 composition but also on its functional properties, as well as the underlying mechanisms.

293 Human faeces constitute in practice the biological samples from which the DNA, RNA  
294 and proteins are extracted in most cases to study the intestinal microbiota composition  
295 and function whereas metabolites and other chemical compounds can be extracted as  
296 well to analyze molecules produced by the microorganisms. Currently, the study of the  
297 IM involves using the new *omics* techniques based on high-throughput sequencing  
298 tools, also called *second-generation sequencing technology*. The DNA sequencing of  
299 the whole IM and the gene functions classifications are performed by *metagenomics*.  
300 *Proteomics* sequence the protein structures to determine cell metabolism through the  
301 activity of the cell enzymes. The analysis of molecules produced by bacterial  
302 metabolism is made by *metabolomics*, and *transcriptomics* studies the complete RNA  
303 molecules quantifying the dynamic expression of genes under different conditions. The  
304 effects of gut microbiota on the host are reflected in different aspects and the  
305 combinations of those *multi-omics* tools provide a new phase in the study of the IM and  
306 its physiological role, linking the composition of the IM with host metabolism, disease  
307 pathogenesis and predictions of therapeutic targets [34].

308

### 309 5. Intestinal microbiota dysbiosis is associated with colorectal cancer and pre- 310 cancerous states

311 Several studies have demonstrated that gut microbiota profiles from CRC patients are  
312 different from that of healthy individuals [35]. Generally, patients with CRC have  
313 decreased microbial diversity in faeces [36] and at the intestinal mucosa level [37]. It is  
314 currently not possible to define a common cancer-associated microbiota [11, 38].  
315 However, although no individual member of the gut microbiota alone is sufficient to  
316 promote CRC, certain microbes have been associated with this type of cancer through  
317 the formation of harmful metabolites and the regulation of certain miRNAs, which then  
318 promote an oncogenic microenvironment. There is evidence of IM associations with  
319 CRC for *Streptococcus bovis*, which has been renamed *Streptococcus gallolyticus*,  
320 *Fusobacterium nucleatum*, *Bacteroides fragilis*, *Enterococcus faecalis* and certain  
321 pathogenic strains from *Escherichia coli* [36]. However, it is not clear at present if these  
322 microorganisms are drivers or passengers in CRC. In addition, although some  
323 microbiota profiles have been associated with the onset and early progression of CRC,  
324 studies in this field are still scarce [39, 40]. Some members of the gut microbiota can  
325 produce microbial genotoxins such as colibactin by *E. coli* group B and fragylisin by *B.*  
326 *fragilis*. Other compounds with cytotoxic action, and potential involvement in the  
327 development of CRC are produced by intestinal microbes such as *Salmonella enterica*,  
328 *Helicobacter pylori*, *F. nucleatum*, *B. fragilis*, *Pseudomonas aeruginosa*,  
329 *Peptostreptococcus anaerobius* and *E. faecalis* among others [11]. The microbial  
330 dysbiosis can also induce changes in host gene expression, subsequently favouring the  
331 development of CRC.

332

### 333 6. Role of the intestinal microbiota on the genotoxic/mutagenic potential of 334 dietary toxic compounds

335 The genotoxicity is the capability to cause damage to the cellular genetic material, and  
336 more specifically mutagenicity is the capacity of genotoxic compounds to alter the  
337 DNA sequence, modifying the expression and functionality of genes. The genotoxicity



338 and/or the mutagenicity in faeces could be determined in an affordable way using some  
339 *in vitro* tests currently available [11].

340 It has been suggested that there is an association of inflammation with the faecal  
341 genotoxicity and CRC through the relationship existing between the gut microbiota and  
342 the innate immune system [38]. Early intestinal mucosal damage (dysplastic lesions,  
343 aberrant crypt foci, and/or intestinal polyps) can precede in years the development of  
344 CRC and these mucosal lesions could be considered early markers of risk for the  
345 development of CRC. Intestinal mucosal lesions are routinely examined for diagnostic  
346 purposes in patients submitted to colonoscopy at hospitals, allowing to differentiate  
347 neoplastic lesions, preneoplastic lesions and healthy intestinal mucosa.

348 The efficiency of endogenous mechanisms of detoxification in the human body largely  
349 depends on the metabolic state of the host, and the type and levels of toxic compounds.  
350 Orally ingested toxic compounds initially reach the liver by direct gut wall absorption  
351 where they are detoxified through phase I (cytochrome P450 system) and phase II  
352 (sulphate, glutathione or glucuronide conjugates) enzymes and are subsequently stored  
353 in the gallbladder. Liver-generated detoxified potential carcinogens are poured again  
354 through the intestine by enterohepatic circulation during digestion (phase III) where  
355 they can be transformed by the gut microbiota.

356 Faecal toxic compounds contributing to genotoxicity may have diverse origins. As  
357 commented before, some members of the intestinal microbiota can produce endogenous  
358 metabolites with genotoxic potential. Other compounds are formed endogenously by the  
359 metabolic activity of intestinal bacteria on dietary constituents such as nitrates, dietary  
360 amines and cholesterol, or are synthesized from precursors of the human metabolism  
361 such as the N-nitroso compounds, fecapentaenes, long-chain fatty acids and secondary  
362 bile acids. The production of these toxic compounds by the IM will depend not only  
363 on the microbiota itself but also on the host physiology, and the interaction of the IM with  
364 diet. In addition, other toxic substances arriving to the gut are of exogenous origin  
365 (foods) and include mycotoxins, plant glycosides, food additives, and the chemical  
366 compounds formed during cooking and food processing commented on previously.

367 Studies using *in vitro* and *in vivo* models indicate that toxic dietary compounds, apart  
368 from their direct effect, could adversely affect the gut microbiota, modifying its  
369 diversity, composition and/or functionality, and affecting host-immunity and  
370 metabolism [35, 41, 42]. The IM can also modify the toxicity of these compounds by i)  
371 decreasing their toxicity through direct binding with the microorganisms and  
372 elimination with faeces, ii) metabolising and transforming them into less toxic  
373 compounds, iii) metabolising and transforming them into more toxically active  
374 molecules, and iv) interfering with detoxifying mechanisms of the host, thus  
375 exacerbating their toxicity [11]. The most notable of these last interactions is that  
376 occurring during enterohepatic circulation when toxic molecules inactivated in phase II  
377 by conjugation to glucuronides in the liver, return to the intestine by enterohepatic  
378 circulation. There, the intestinal microbial glucuronidases, mostly from Enterobacteria,  
379 *Clostridium* and *Bacteroides* members, release the inactivated chemical compound from  
380 the glucuronide and subsequently turn it back into a toxic molecule.

381 Global diet modulates the composition and functionality of the IM, influencing the way  
382 in which this microbial community interacts with dietary toxic compounds and with  
383 detoxifying mechanisms of the host, then contributing to increase or decrease in the  
384 intestinal toxicity. In this scenario, it would be possible to identify early shifts in  
385 microbiota patterns (composition and/or functionality) associated at variable degree

386 with increased **intestinal** toxicity, **the** intake **of chemicals with carcinogenic potential**  
387 and global diet. These modifications of the microbiota (even when they could represent  
388 adaptive processes) may **be associated with abnormal changes of the intestinal mucosa**  
389 **that would represent an augmented risk for the subsequent development of CRC. The**  
390 **diversity of chemical structures of dietary toxic compounds and the difficulty to**  
391 **determine accurately their intake with diet substantially increase the challenge of**  
392 **teasing out individual chemical class influences on CRC. However, initial effort like**  
393 **those focusing on a specific and defined group of compounds, as those chemicals**  
394 **generated during food processing, would make the task more realistic and affordable.**  
395 **These compounds could be assessed by means of dietary interviews that include**  
396 **cooking/preparation procedures, duration and temperature of the process, and the use of**  
397 **specific food composition databases.**

398 **Our hypothesis is that beyond differences in genetic susceptibilities, metabolic states**  
399 **and the inherent variability of microbiota profiles among individuals and human groups,**  
400 **the net exposure to dietary molecules with carcinogenic potential will depend on the**  
401 **type of compound, doses, frequency of consumption and lifetime exposure. These**  
402 **factors will be modified by food preparation procedures, which will be closely related**  
403 **to the amount of compound ingested, the global dietary patterns and IM profile of**  
404 **subjects. Therefore, risk thresholds for CRC could be established as a function of gut**  
405 **genotoxicity, IM and diet (global dietary patterns and toxic molecules intake),**  
406 **considering precancerous or cancerous mucosal changes as an outcome variable.**

407 ML and Semantic Web are important tools that could assist in the treatment and  
408 modelling of such data in order to categorize the risk (Fig. 1).

409 The identification of changes in the microbiota associated **with the intake of toxic**  
410 **compounds with carcinogenic potential** could be useful to elaborate guidelines **for** food  
411 processing and dietary recommendations.

412

## 413 **7. ML: a tool to assess risk by dietary exposure**

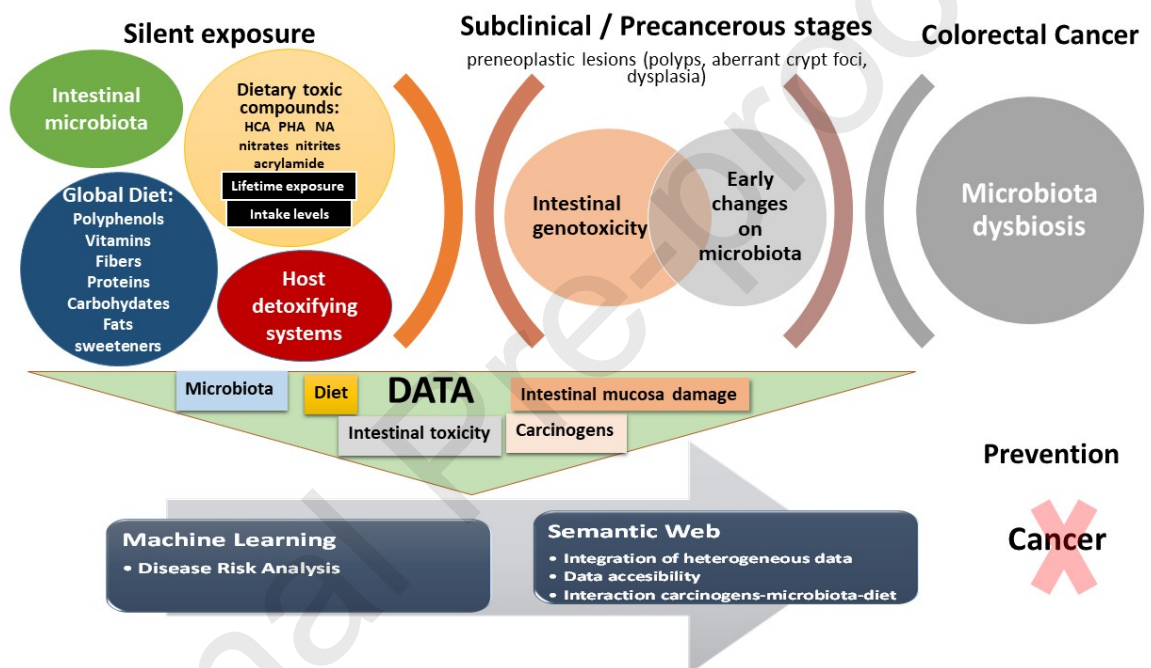
414 ML can be considered a branch of artificial intelligence, **as it attempts** to use computers  
415 to complement human intelligence [43]. ML **has** become an essential tool for  
416 biomedical research and **the** modern healthcare system, given that the amount of  
417 medical and biological data requiring analysis has increased abruptly in **the** last years,  
418 and some ML methods have shown their ability for solving complex problems.

419 A key objective of any learning algorithm is to build models with good generalization  
420 capability [44]. Thus, the classification procedure is a cornerstone in any predictive  
421 problem. In addition, there is not a standard classification method **to date**. Different  
422 methods could be applied to design the prediction model. A decision tree (DT) is a  
423 mathematical tree where the internal nodes are tests on the variables that define the  
424 inputs and the leaf nodes are classes. C5.0, C4.5, CART or Random Forests (RF) are  
425 examples of this kind of ML. Lazy learners such as k-Nearest Neighbours (KNN) are  
426 based on learning by comparing a given test example with each training example.  
427 Artificial Neural Networks (ANN) are inspired in biological neural networks. Kernel  
428 methods as Support Vector Machines (SVM) are based on the idea of embedding the  
429 data into a high dimensional feature space using the kernel [45].

430 ML has been applied to dietary studies and for deciphering the effect of the exposure to  
431 pollutants and carcinogens. Thus, Chatterjee et al. [46] identified potential risk factors  
432 for preventing obesity using a broad set of different ML techniques. In another work

433 [47] the mutual interactions between diet, microbiota, metabolic responses and the  
 434 immune system were developed using a ML-based method. In a similar way, we  
 435 employed DT to study the interactions between serum free fatty acids and faecal  
 436 microbiota [48]. Gut microbiota was also identified as a factor in predicting  
 437 personalised postprandial glycaemic response to real-life meals, obtaining an accurate  
 438 prediction with boosting DT [49]. An oral malodour classifier was developed as a  
 439 function of the oral microbiota in saliva, with SVM, ANN and DT, and SVM being the  
 440 most accurate [50]. The decline of *Akkermansia muciniphila* was identified as a  
 441 common dysbiotic marker linked to disease status by using DTs [51]. Cammarota et al.  
 442 [52] recently highlighted the importance of the gut microbiome and the need of  
 443 applying ML to analyse the considerably quantity of complex health care data in cancer  
 444 research.

445



446

447

448 **Figure 1. Schematic representation of risk assessment by exposure to dietary**  
 449 **toxic compounds formed during food cooking and processing as a function of**  
 450 **the IM, diet and intestinal toxicity, applying ML and Semantic Web.** The net  
 451 exposure to toxic compounds depends on the intake and time of exposure and this  
 452 influences the genotoxicity at the intestinal environment. IM and global diet could  
 453 modify the resulting toxicity of dietary chemicals. Prolonged exposure to high  
 454 intestinal toxicity levels could lead to changes in the intestinal mucosa that may be  
 455 accompanied by shifts in the intestinal microbiota. Applying ML to dietary and  
 456 microbiota data in silent, subclinical and precancerous stages of intestinal mucosal  
 457 damage could assist in CRC risk assessment whereas Semantic Web will facilitate  
 458 data accessibility and management.

459 Therefore, ML has proven to be an efficient tool to identify some key factor  
 460 relationships associated with diet, health parameters and lifestyles with the microbiota  
 461 and disease [48-51]. Although no general rule exists *a priori* indicating which ML  
 462 method is the best, depending on a given problem, it is expected that ML could  
 463 successfully contribute to establishing risk thresholds for CRC as a function of the

464 intake of chemicals with carcinogenic potential, global diet, intestinal genotoxicity and  
 465 shifts in microbiota profiles. In summary, ML is able to consider factors from different  
 466 sources (such as those related to ingested of potential carcinogens, diet and IM), select  
 467 the most relevant ones and use them to predict the risk of CRC. A general workflow of  
 468 the process is provided in Fig. 2.

469

## 470 8. Worked example of a ML process for CRC risk assessment

471 Since real data on diet, intake of toxic chemicals, intestinal microbiota and fecal  
 472 genotoxicity/mutagenicity are not yet available in a single database, a conceptual design  
 473 is proposed using previously published variables corresponding to the metabolism of  
 474 healthy people and people with CRC.

475 *Dataset.* The dataset employed is a subset of the Colorectal Cancer Detection Using  
 476 Targeted Serum Metabolic Profiling experiment from University of Washington. These  
 477 data are available at <https://www.metabolomicsworkbench.org/>.

478 The dataset is composed by 234 individuals and 124 variables. For this example,  
 479 Diagnosis is the target variable, that is recoded as a binary variable representing if each  
 480 example presents colorectal cancer or not. Since real data are not yet available, a  
 481 conceptual design is proposed using previously published variables corresponding to the  
 482 metabolism of healthy people and CRC. From the total of existing variables in the  
 483 repository, we have selected those that could be directly correlated with the diet (sugars,  
 484 aminoacids, fatty acids and other compounds of interest) and including some  
 485 anthropometrical variables related with diet and health, as the BMI. In addition, from  
 486 the 124 variables, we have selected the following as predictive ones to run this example:  
 487 "Acetylcholine" "Alanine" "Asparagine" "Aspartic\_Acid" "Biotin" "Glutamic\_acid"  
 488 "Glutamine" "Histidine" "Linolenic\_Acid" "Lysine" "Methionine" "Methylsuccinate"  
 489 "Pyruvate" "Tryptophan" "BMI"

490 The following tables show basic statistics for these variable set depending on the value  
 491 of the target variable.

VARIABLE	HEALTHY		
	min	mean	max
ACETYLCHOLINE	227140.38	1944056.93	3933866.8
ALANINE	4029094.49	6339425.03	10736506.9
ASPARAGINE	446544.10	697142.28	926673.2
ASPARTIC_ACID	367199.83	1207280.26	2736972.2
BIOTIN	70262.10	134817.68	218108.1
GLUTAMIC_ACID	696905.02	2101133.59	4333471.5
GLUTAMINE	23520246.16	32222162.15	42570355.8
HISTIDINE	10280560.84	18694498.58	29272649.1
LINOLENIC_ACID	403397.25	865422.94	1610396.2
LYSINE	5435894.46	10117619.23	13735189.6
METHIONINE	306713.00	732652.67	1004676.9
METHYLSUCCINATE	801214.57	1303371.84	1856837.4
PYRUVATE	55107.82	174507.45	429810.3
TRYPTOPHAN	501607.00	3715594.49	5471963.8

**BMI** | 20.00 | 27.58 | 42.0

492

<b>COLORECTAL CANCER</b>			
<b>VARIABLE</b>	<b>min</b>	<b>mean</b>	<b>max</b>
<b>ACETYLCHOLINE</b>	712642.00	1755303.59	3723973.0
<b>ALANINE</b>	2910976.71	5640811.77	9555174.6
<b>ASPARAGINE</b>	456356.62	656879.58	1052985.7
<b>ASPARTIC_ACID</b>	377375.66	1636515.34	4411499.3
<b>BIOTIN</b>	63989.62	123128.62	228928.5
<b>GLUTAMIC_ACID</b>	916836.31	2683576.11	6559485.0
<b>GLUTAMINE</b>	16182419.38	29168842.58	36269190.5
<b>HISTIDINE</b>	8189632.57	14905491.46	25936858.0
<b>LINOLENIC_ACID</b>	167055.75	662328.07	1213540.7
<b>LYSINE</b>	5237148.72	8703904.55	12749510.4
<b>METHIONINE</b>	338104.80	617976.05	1045772.3
<b>METHYLSUCCINATE</b>	825623.96	1207703.72	1885528.9
<b>PYRUVATE</b>	64219.18	199196.83	458775.4
<b>TRYPTOPHAN</b>	1785060.16	3451357.71	5410601.2
<b>BMI</b>	17.00	25.35	32.0

493

494 *Preprocessing.* As it is well known that some ML methods are quite sensitive to  
 495 variable scale, continuous variables were normalized. In addition, missing values were  
 496 treated using K-nearest neighbor imputation.

497 *Classification and evaluation.* As it was highlighted before, a key objective of any  
 498 learning algorithm is to build models with good generalization capability, which is  
 499 equivalent to look for models that accurately predict the class labels of previously  
 500 unknown examples. Therefore, the classification procedure is a cornerstone in any  
 501 predictive problem. In addition, there is no a standard classification method so far.  
 502 Thus, several different methods were tested to select the one performing the best for this  
 503 task, taking into account the trade-off between performance and interpretability. The  
 504 methods considered in this worked example are a tree based method (C4.5), a lazy  
 505 learners (Knn), a Neural Network (in particular, multilayer perceptrons, MLP) and a  
 506 support vector machine with radial kernel.

507 Training a ML method is as complex as necessary to avoid overfitting and to correctly  
 508 optimize the different hyperparameters associated to each method. In this case we have  
 509 applied cross-validation with 10 folds. During the cross validation process, the specific  
 510 parameters associated to each method have been optimized using the default  
 511 configuration.

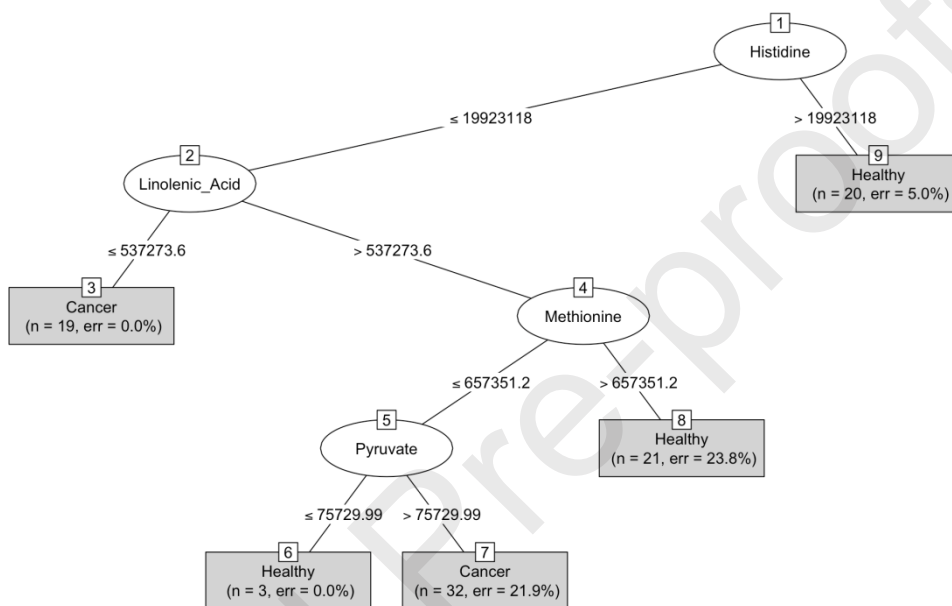
512

513

	<b>Sensitivity</b>	<b>Specificity</b>
<b>J48</b>	0.75	0.63
<b>SvmRadial</b>	0.80	0.62
<b>Knn</b>	0.87	0.87
<b>MLP</b>	0.71	0.64

514

515 From the results obtained, it is clear that the method performing better according to both  
 516 Sensitivity and Specificity is KNN. The value of k was 9. Note that this parameter is set  
 517 experimentally in training phase. It is well known that KNN does not provide  
 518 information about the features providing this classification. Thus, using this method, it  
 519 is only possible to predict if an example is labelled as Healthy or having CRC. The  
 520 same occurs with MLP and SvmRadial. As a consequence, if one is interested in  
 521 analyzing the factors helping in the prediction, a model based on decision trees should  
 522 be selected. The one employed here is C4.5. In this example, the model produced is the  
 523 following:



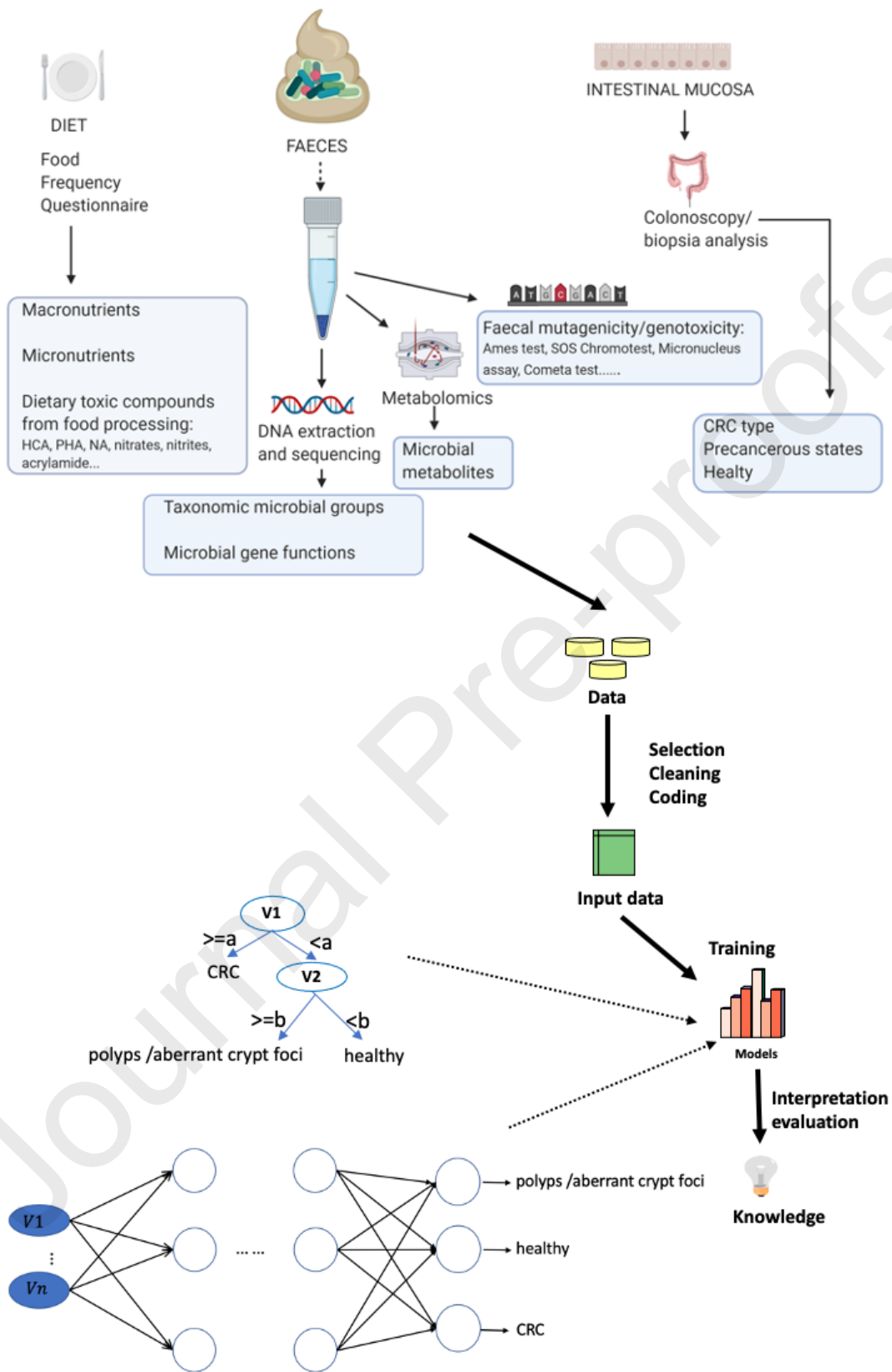
524

525 From the initial set of variables, "Acetylcholine", "Alanine", "Asparagine",  
 526 "Aspartic\_Acid", "Biotin", "Glutamic\_acid", "Glutamine", "Histidine",  
 527 "Linolenic\_Acid", "Lysine", "Methionine", "MethylSuccinate", "Pyruvate",  
 528 "Tryptophan", "BMI", C4.5 detects Histidine, Linolenic\_Acid, Methionine and Pyruvate  
 529 as relevant variables for predicting CRC.

530 All the experiments in this worked example were performed using RStudio 1.3.1093, R  
 531 4.0.3 and caret package, version 6.0-86.

532

533 **Figure 2. General workflow of a Machine Learning process for CRC risk**  
 534 **assessment as a function of diet, microbiota and intestinal genotoxicity.** Data  
 535 from diet (FFQ), microbial metabolites, microbiota composition, microbial gene  
 536 functions, and genotoxicity/mutagenicity (faeces) and biopsy analyses of the  
 537 intestinal mucosa (routine colonoscopies at hospitals) are collected in a joint  
 538 database and submitted to a ML process. Some ML models (such as DT, on bottom-  
 539 left) allow establishing profiles and thresholds related to the input variables, while  
 540 others (such as ANN, on bottom-right) are more difficult to interpret but are  
 541 successful predictors.

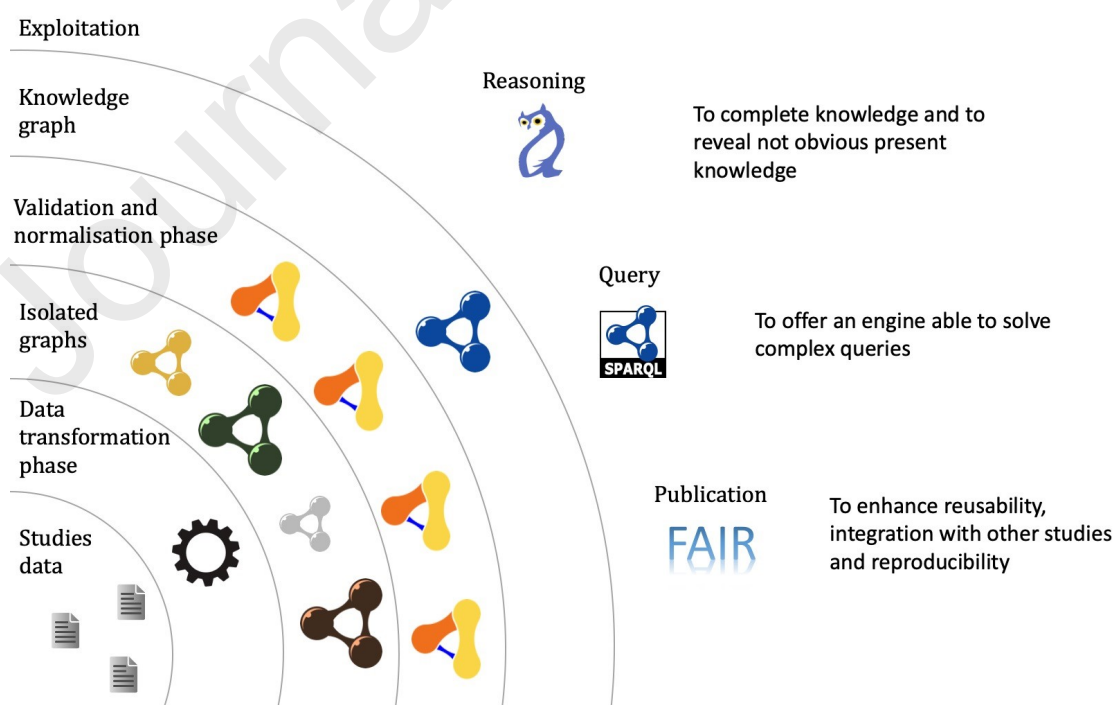


## 543 9. Using Semantic Web to connect and to exploit data

544 The Semantic Web vision has supposed a shift of persistence, modelling and  
 545 interoperability of data [53]. Being able to represent entities unambiguously, link them  
 546 and integrate different data-sources in a single representation, has enabled a new set of  
 547 semantic-aware applications. These computer science advances are ready to be applied  
 548 to different fields. Specifically, in the bio-computational field, some works have  
 549 explored its use i) to describe human and mouse genes [54], ii) to offer a platform that  
 550 eases the consumption and curation of genome data [55], iii) to integrate different drug  
 551 data-sources [56], iv) to provide a platform to analyse the course of diseases [57].  
 552 Therefore, we envisage next challenges using Semantic Web technologies to model and  
 553 to exploit data from nutrition and microbiota interaction studies (Fig.3).

554 One of the main problems facing the exploitation of data from these type of studies is  
 555 the existence of many heterogeneous data-sources with their own data models that  
 556 cannot be integrated easily with others. This issue prevents obtaining conclusions of the  
 557 joint-analysis of data from different studies. To alleviate this problem, some ontologies  
 558 were proposed which ensure that all data providers are talking about the same domain  
 559 [58]. For example, FoodOn [59] for data integration of food traceability and quality  
 560 control is a very specific ontology that offers a great basis for reusability. In contrast,  
 561 ONS [60] is a general ontology for nutrition studies that can be tuned with specific  
 562 elements if necessary. Alongside the creation of well-defined ontologies, there arises the  
 563 need for tools able to migrate non-semantic data to these new semantic standards.  
 564 Recent development of heterogeneous data mapping tools in the Semantic Web has  
 565 supposed a new paradigm in knowledge graph creation methodologies [61], offering  
 566 reusability, maintainability and a better user-experience. The use of these tools can  
 567 deliver a faster migration of non-semantic datasets to a knowledge graph in which all  
 568 desired studies can be integrated. This will offer the possibility to analyse all data  
 569 together, make it accessible, and preserve it for future uses, which is in keeping with  
 570 FAIR (Findable, Accessible, Interoperable and Reusable data) principles [62].

571



572



**Figure 3. Semantic Web schema and technological stack proposed for microbiota and diet studies.** Each concentric circumference represents a layer/process in the technological stack; these layers are independent and can work by themselves. The layer stacking means that an upper layer contains the lower ones and need for them to be complete and coherent. Different coloured graphs represent graphs from different sources, which are not yet integrated. Orange and yellow patterns in the validation phase represent the mechanism of validation and normalization of the aforementioned heterogeneous graphs, which connect to a unique and integrated knowledge graph.

Although a well-defined ontology can enable interoperability and integration of different datasets, we must also ensure that different pieces of data follow the same shape, which will derive in a cleaned and normalised graph and, therefore, an easier one to query. The use of Resource Description Framework (RDF) [63] validation technologies was explored in Fast Healthcare Interoperability Resources (FHIR) specification [64] to not only validate data but to share data models among humans and machines [65]. Therefore, using ontologies, we can define the meta-knowledge of the domain, e.g., the category's relationships between different mutagens, nutrients or bacteria; using RDF validation techniques we can ensure certain rules, e.g., that a value is between certain limits or that a nutrient has a certain number of attributes.

Once various datasets are converted, validated—using the aforementioned techniques—and their semantics defined using a proper ontology, new results could be delivered. Thanks to ontology axioms it is possible to generate inferences on pre-existing knowledge in order to reveal non-evident and underlying content, which could be obviated [66]. For example, if we define *Bacteroides fragilis* we know that it also belongs to the categories *Bacteroides* (genus), *Bacteroidaceae* (family), *Bacteroidales* (order), *Bacteroidia* (class) and *Bacteroidetes* (phylum); however, this information is not evident for a machine. Thus, the inference system will fill these upper categories, so all data is complete and can be easily integrated. In addition, the graph data model used by RDF enables a different data modelling—in contrast with the normally used tabular form—that by means of SPARQL—the advocated RDF query language—could reveal new relationships previously obviated [67]. This simplifies the modelling of the former example in which we have multiple categories, and consequently we wish that *B. fragilis* were shown when asking for a *Bacteroidetes*, and a *Bacteroidaceae*, among others. Doing the same modelling in tabular form would imply considerably more complicated structures that can be error-prone.

Finally, this methodology offers the possibility to not only improve analysis techniques and discover hidden content but also to transfer part of this knowledge and make it accessible for the public. The emergence of projects as Wikidata [68] enables the creation of general-purpose knowledge graphs integrating data that could be interesting for the entire world and that is curated by users. It is possible, by taking advantage of proposed conversions, to publish interesting conclusions of involved studies in the so-called semantic eScience [69]. This approach may be employed for the achievement of FAIR principles but also to achieve a transference and dissemination effort, which could lead to a relief in the ongoing reproducibility crisis [70].

## 620 10. Summary and Outlook

621 The net exposure to dietary **toxic compounds**, and the intestinal genotoxicity generated,  
622 depends on the intake and time of consumption and on their interaction with the IM and  
623 global diet. The IM of individuals with CRC differs from that of healthy people, but  
624 studies relating the consumption of carcinogens with adverse early shifts of microbiota  
625 (either beneficial adaptive or adverse changes) are very scarce. The complexity of data  
626 and the **several** variables potentially affecting these interactions may hinder the  
627 interpretation of the studies. In this context, the application of ML to the data obtained  
628 in subclinical and precancerous stages **of the intestinal mucosa** could help to analyse the  
629 risk **for development of CRC associated to** the intake of carcinogens as a function of  
630 diet and microbiota profiles. Moreover, the use of the recently developed Semantic Web  
631 approaches could improve data accessibility and management, contributing to evidence  
632 **of** new interactions among carcinogens, microbiota, and diet (Fig. 1).

633

### 634 **Competing financial interest**

635 The authors declare no competing financial interest or personal relationships that could  
636 have influenced the content of this article.

637

### 638 **CRedit authorship contribution statement**

639 **Sergio Ruiz Saavedra:** Writing-original draft, review & editing. **Herminio García**  
640 **González:** Writing-original draft, review & editing. **Silvia Arboleya:** Writing-original  
641 draft, review & editing. **Nuria Salazar:** Writing-original draft, review & editing. **Jose**  
642 **Emilio Labra-Gayo:** Writing-original draft, review & editing. **Susana Irene Díaz:**  
643 Writing-original draft, review & editing. **Miguel Gueimonde:** Writing-original draft,  
644 review & editing. **Sonia González:** Conceptualization, Funding acquisition, Writing-  
645 original draft, review & editing. **Clara G. de los Reyes-Gavilán:** Supervision,  
646 Conceptualization, Funding acquisition, Writing-original draft, review & editing

647

648

### 649 **Acknowledgements**

650 This work is receiving support from the project RTI2018-098288-B-I00  
651 (MCIU/AEI/FEDER, UE) and is based on concepts developed partly funded by projects  
652 TIN2017-88877-R (AEI/FEDER, UE), TIN2017-87600-P (AEI/FEDER, UE) and  
653 IDI/2018/000176 (Asturian Government GRUPIN projects).

654 SRS is the beneficiary of a training contract financed by project RTI2018-098288-B-  
655 I00. SA is granted by a Juan de la Cierva postdoctoral contract from the Spanish  
656 Ministry of Science and Innovation (Ref. IJCI-2017-32156) and NS is the recipient of a  
657 postdoctoral contract awarded by the Fundación para la Investigación y la Innovación  
658 Biosanitaria del Principado de Asturias (FINBA).

659 Figure 2 was partly created with Biorender.com

660

661

662

663

664 **References:**

- 665 [1] Bouvard V, Loomis D, Guyton KZ, Grosse, Y, El Ghissassi F, et al. (2015)  
666 Carcinogenicity of consumption of red and processed meat. *Lancet Oncol* 16: 1599–  
667 600.
- 668 [2] Sugimura T, Wakabayashi K, Nakagama H, Nagao M (2004) Heterocyclic amines:  
669 mutagens/carcinogens produced during cooking of meat and fish. *Cancer Sci* 95:  
670 290–9.
- 671 [3] Chiavarini M, Bertarelli G, Minelli L, Fabiani R (2017) Dietary intake of meat  
672 cooking-related mutagens (HCAs) and risk of colorectal adenoma and cancer: A  
673 systematic review and meta-analysis. *Nutrients* 9: 514-36.
- 674 [4] Zheng W, Lee SA (2009) Well-done meat intake, heterocyclic amine exposure, and  
675 cancer risk. *Nutr Cancer* 61: 437–46.
- 676 [5] Gibis M (2016) Heterocyclic aromatic amines in cooked meat products: causes,  
677 formation, occurrence, and risk assessment. *Compr Rev Food Sci Food Saf* 15: 269–  
678 302.
- 679 [6] Zimmerli B, Rhyn P, Zoller O, Schlatter J (2001) Occurrence of heterocyclic  
680 aromatic amines in the Swiss diet: Analytical method, exposure estimation and risk  
681 assessment. *Food Addit Contam* 18:533–51.
- 682 [7] Ericson U, Wirfält E, Mattisson I, Gullberg B, Skog K (2007) Dietary intake of  
683 heterocyclic amines in relation to socio-economic, lifestyle and other dietary factors:  
684 Estimates in a Swedish population. *Public Health Nutr* 10: 616–27.
- 685 [8] Butler LM, Sinha R, Millikan RC, Martin CF, Newman B, et al. (2003)  
686 Heterocyclic amines, meat intake, and association with colon cancer in a population-  
687 based study. *Am J Epidemiol* 157: 434–45.
- 688 [9] Wong KY, Su J, Knize MG, Koh WP, Seow A (2005) Dietary exposure to  
689 heterocyclic amines in a Chinese population. *Nutr Cancer* 52: 147–55.
- 690 [10] Pouzou, JG, Costard, S, Zagmutt, FJ (2018) Probabilistic estimates of heterocyclic  
691 amines and polycyclic aromatic hydrocarbons concentrations in meats and breads  
692 applicable to exposure assessments. *Food Chem Toxicol* 114: 346–60.
- 693 [11] Nogacka AM, Gómez-Martín M, Suárez A, González-Bernardo O, de los Reyes-  
694 Gavilán CG, et al. (2019) Xenobiotics formed during food processing: their relation  
695 with the intestinal microbiota and colorectal cancer. *Int J Mol Sci* 20: 2051.
- 696 [12] Yu YX, Chen L, Yang D, Pang YP, Zhang SH, et al. Polycyclic aromatic  
697 hydrocarbons in animal-based foods from Shanghai: Bioaccessibility and dietary  
698 exposure. *Food Addit Contam - Part A Chem Anal Control Expo Risk Assess*  
699 29:1465–74.
- 700 [13] Domingo JL, Nadal M (2015) Human dietary exposure to polycyclic aromatic  
701 hydrocarbons: A review of the scientific literature. *Food Chem Toxicol* 86: 144–53.
- 702 [14] Alexander J, Benford D, Cockburn A, Cravedi J P, Dogliotti E, et al (2008)  
703 Polycyclic Aromatic Hydrocarbons in food - Scientific opinion of the panel on  
704 contaminants in the food chain. *EFSA J* 724: 1-114.
- 705 [15] IARC Working Group on the Evaluation of Carcinogenic Risks to Humans (2012)  
706 Chemical agents and related occupations. *IARC Monogr Eval Carcinog Risks Hum*  
707 100: 9–562.
- 708 [16] Steinberg P (2019) Red meat-derived nitroso compounds, lipid peroxidation  
709 products and colorectal cancer. *Foods* 8: 252.
- 710 [17] Zhang FX, Miao Y, Ruan JG, Meng SP, Dong J Da, et al. (2019) Association  
711 between nitrite and nitrate intake and risk of gastric cancer: A systematic review and  
712 meta-analysis. *Med Sci Monit* 25: 1788–99.

- 713 [18] Zheng J, Stuff J, Tang H, Hassan MM, Daniel CR, et al (2019) Dietary N-nitroso  
714 compounds and risk of pancreatic cancer: results from a large case-control study.  
715 *Carcinogenesis* 40: 254–262.
- 716 [19] Song P, Wu, L, Guan W (2015) Dietary nitrates, nitrites, and nitrosamines intake  
717 and the risk of gastric cancer: a meta-analysis. *Nutrients* 7: 9872-9895.
- 718 [20] Zhu Y, Wang PP, Zhao J, Green R, Sun Z, et al (2014) Dietary N-nitroso  
719 compounds and risk of colorectal cancer: a case-control study in Newfoundland  
720 and Labrador and Ontario, Canada. *Br J Nutr* 111: 1109-1117.
- 721 [21] Anese M, Nicoli M C, Verardo G, Munari M, Mirolo G (2014) Effect of vacuum  
722 roasting on acrylamide formation and reduction in coffee beans. *Food Chem* 145:  
723 168–172.
- 724 [22] European Food Safety Authority (2011) Results on acrylamide levels in food from  
725 monitoring years 2007–2009 and exposure assessment. *EFSA J* 9: 2133.
- 726 [23] EFSA Panel on Contaminants in the Food Chain (CONTAM) (2015) Scientific  
727 opinion on acrylamide in food. *EFSA J* 13: 4104.
- 728 [24] Cuparencu C, Praticó G, Hemeryck LY, Sri Harsha PSC, Noerman S, et al. (2019)  
729 Biomarkers of meat and seafood intake: An extensive literature review. *Genes Nutr*  
730 14: 1–30.
- 731 [25] Jakszyn P, Ibáñez R, Pera G, Agudo A, García-Closas R, et al. (2004) Food content  
732 of potential carcinogens. Barcelona: Catalan Institute of Oncology. Available from:  
733 <http://epic-spain.com/libro.html>.
- 734 [26] National Cancer Institute (2006). CHARRED: computerized heterocyclic amines  
735 database resource for research in the epidemiologic of disease. Available from  
736 <https://dceg.cancer.gov/tools/design/charred> (Accessed: 20 July 2020).
- 737 [27] European Food Safety Authority (2012). Update on acrylamide levels in food from  
738 monitoring years 2007 to 2010. *EFSA Journal* 10(10):2938.
- 739 [28] Thursby E, Juge N (2017) Introduction to the human gut microbiota. *Biochem J*  
740 474: 1823–36.
- 741 [29] Ríos-Covián D, Ruas-Madiedo P, Margolles A, Gueimonde M, De los Reyes-  
742 Gavilán C G, et al. (2016) Intestinal short chain fatty acids and their link with diet  
743 and human health. *Front Microbiol* 7: 185.
- 744 [30] Gao J, Xu K, Liu H, Liu G, Bai M, et al. (2018) Impact of the gut microbiota on  
745 intestinal immunity mediated by tryptophan metabolism. *Front Cell Infect Microbiol*  
746 8: 13.
- 747 [31] Rios-Covian D, González S, Nogacka A M, Arboleya S, Salazar N, et al. (2020) An  
748 Overview on fecal branched short-chain fatty acids along human life and as related  
749 with body mass index: Associated dietary and anthropometric factors. *Front*  
750 *Microbiol* 11: 973.
- 751 [32] Turnbaugh PJ, Hamady M, Yatsunenko T, Cantarel BL, Duncan A, et al. (2009) A  
752 core gut microbiome in obese and lean twins. *Nature* 457: 480–4.
- 753 [33] Echarri PP, Graciá CM, Berruezo GR, Vives I, Ballesta M, et al. (2011).  
754 Assessment of intestinal microbiota of full-term breast-fed infants from two  
755 different geographical locations. *Early Hum Dev* 87: 511-13 .
- 756 [34] Wang X, Zhang A, Miao J, Sun H, Yan G, et al. (2018) Gut microbiota as important  
757 modulator of metabolism in health and disease. *RSC Advances*, 8: 4239.
- 758 [35] Abu-Ghazaleh N, Chua WJ, Gopalan V (2020) Intestinal microbiota and its  
759 association with colon cancer and red/processed meat consumption. *J Gastroenterol*  
760 *Hepatol*.
- 761 [36] Ahn J, Sinha R, Pei Z, Dominianni C, Wu J, et al. (2013) Human gut microbiome  
762 and risk for colorectal cancer. *J Natl Cancer Inst* 105: 1907–11

- 763 [37] Chen W, Liu F, Ling Z, Tong X, Xiang C (2012) Human intestinal lumen and  
764 mucosa-associated microbiota in patients with colorectal cancer. PLOS ONE 7:  
765 e39743.
- 766 [38] Maissonneuve C, Irrazabal T, Martin A, Girardin SE, Philpott DJ (2018) The impact  
767 of the gut microbiome on colorectal cancer. Annu Rev Cancer Biol 2: 229-49.
- 768 [39] Kinross J, Mirnezami R, Alexander J, Brown R, Scott A, et al. (2017) A prospective  
769 analysis of mucosal microbiome-metabonome interactions in colorectal cancer using  
770 a combined MAS 1HNMR and metataxonomic strategy. Sci Rep 7: 8979.
- 771 [40] Pu LZ, Yamamoto K, Honda T, Nakamura M, Yamamura T, et al. (2020)  
772 Microbiota profile is different for early and invasive colorectal cancer and is  
773 consistent throughout the colon. J Gastroenterol Hepatol 35: 433–7.
- 774 [41] Beer F, Urbat F, Franz CMAP, Huch M, Kulling SE, et al. (2019) The human fecal  
775 microbiota metabolizes foodborne heterocyclic aromatic amines by reuterin  
776 conjugation and further transformations. Mol Nutr Food Res 63: 1801177.
- 777 [42] Ribière C., Peyret P, Parisot N, Darcha, C, Déchelotte PJ, et al. (2016) Oral  
778 exposure to environmental pollutant benzo[a]pyrene impacts the intestinal  
779 epithelium and induces gut microbial shifts in murine model. Sci Rep 6: 31027.
- 780 [43] Shalev-Shwartz S, Ben-David S (2014) Understanding Machine Learning: From  
781 theory to algorithms, Cambridge University Press.
- 782 [44] Mitchell T (1997) Machine Learning, McGraw-Hill.
- 783 [45] Kuhn M, Johnson K (2013) Applied predictive modeling, Springer
- 784 [46] Chatterjee A, Gerdes MW, Martinez SG (2020) Identification of risk factors  
785 associated with obesity and overweight-a machine learning overview. Sensors:  
786 2734.
- 787 [47] Danneskiold-Samsøe NB, Dias de Freitas H, Santos R, Lemos Bicas J, Cazarin C, et  
788 al. (2019) Interplay between food and gut microbiota in health and disease. Food  
789 Res Int 115: 23-31.
- 790 [48] Fernandez-Navarro T, Díaz I, Gutiérrez-Díaz I, Rodríguez-Carrio J, Suárez A, et al.  
791 (2019) Exploring the interactions between serum free fatty acids and fecal  
792 microbiota in obesity through a machine learning algorithm. Food Res Int 121: 533–  
793 541.
- 794 [49] Zeevi D, Korem T, Zmora N, Israeli D, Rothschild D, et al. (2015) Personalized  
795 nutrition by prediction of glycemic responses. Cell 163: 1079–1094.
- 796 [50] Nakano Y, Takeshita T, Kamio N, Shiota S, Shibata Y, et al. (2014) Supervised  
797 machine learning-based classification of oral malodor based on the microbiota in  
798 saliva samples. Artif Intell Med 60: 97–101.
- 799 [51] Lopetuso LR, Quagliariello A, Schiavoni M, Petito V, Russo A, et al. (2020)  
800 Towards a disease-associated common trait of gut microbiota dysbiosis: The pivotal  
801 role of *Akkermansia muciniphila*. Digest Liver Dis 52: 1002-1010.
- 802 [52] Cammarota G, Ianiro G, Ahern A, Carbone C, Temko A, et al. (2020) Gut  
803 microbiome, big data and machine learning to promote precision medicine for  
804 cancer. Nat Rev Gastroenterol Hepatol 17:635-648.
- 805 [53] Berners-Lee T, Hendler J, Lassila O (2001) The semantic Web. Scientific American  
806 284:34-43.
- 807 [54] Burgstaller-Muehlbacher S, Waagmeester A, Mitraka E, Turner J, Putman T, et al.  
808 (2016) Wikidata as a semantic framework for the gene wiki initiative. Database Vol  
809 2016: baw015.
- 810 [55] Putman TE, Lelong S, Burgstaller-Muehlbacher S, Waagmeester A, Diesh C, et al.  
811 (2017) WikiGenomes: an open web application for community consumption and  
812 curation of gene annotation data in Wikidata. Database Vol 2017: bax025.

- 813 [56] Gray AJ, Askjaer S, Brenninkmeijer CY, Burger K, Chichester C, et al. (2012) The  
814 Pharmacology Workspace: A platform for drug discovery. Proceedings of the 3rd  
815 International Conference on Biomedical Ontology (ICBO 2012), KR-MED Series,  
816 Graz, Austria, July 21-25, 2012; CEUR Workshop Proceedings; 897.
- 817 [57] Esteban-Gil A, Fernández-Breis JT, Boeker M (2017) Analysis and visualization of  
818 disease courses in a semantically-enabled cancer registry. *J Biomed Semantics*. 8:  
819 46.
- 820 [58] Chandrasekaran B, Josephson JR, Benjamins VR (1999) What are ontologies, and  
821 why do we need them?. *IEEE Intell Syst App* 14: 20-26.
- 822 [59] Dooley DM, Griffiths EJ, Gosal GS, Buttigieg PL, Hoehndorf R, et al. (2018)  
823 FoodOn: a harmonized food ontology to increase global food traceability, quality  
824 control and data integration. *NPJ SciFood* 2: 23.
- 825 [60] Vitali F, Lombardo R, Rivero D, Mattivi F, Franceschi P, et al. (2018) ONS: an  
826 ontology for a standardized description of interventions and observational studies in  
827 nutrition. *Genes Nutr* 13: 12.
- 828 [61] De Meester B, Heyvaert P, Verborgh R, Dimou A (2019) Mapping languages  
829 analysis of comparative characteristics. First Knowledge Graph Building Workshop,  
830 part of ESWC2019. Portorož, Slovenia, June 3, 2019; CEUR Workshop  
831 Proceedings 2489: 37-45.
- 832 [62] Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, et al. (2016)  
833 The FAIR Guiding Principles for scientific data management and stewardship.  
834 *Scientific data* 3: 160018.
- 835 [63] Cyganiak R, Wood D, Lanthaler M (2014) RDF 1.1 Concepts and Abstract Syntax.  
836 W3C Recommendation. 25 February 2014. Available online:  
837 <https://www.w3.org/TR/rdf11-concepts/>.
- 838 [64] Bender D, Sartipi K (2013) HL7 FHIR: An Agile and RESTful approach to  
839 healthcare information exchange. Proceedings of the 26th IEEE international  
840 symposium on computer-based medical systems. Porto, Portugal, June 20-22, 2013.  
841 IEEE Computer Society 326-31.
- 842 [65] Thornton K, Solbrig H, Stupp GS, Labra Gayo JE, Mietchen D, et al. (2019) Using  
843 Shape Expressions (ShEx) to share RDF data models and to guide curation with  
844 rigorous validation. The Semantic Web - 16th International Conference, ESWC  
845 2019, Portorož, Slovenia, June 2-6, 2019, Springer Lecture Notes in Computer  
846 Science 11503: 606-20.
- 847 [66] Khamparia A, Pandey B (2017) Comprehensive analysis of semantic web reasoners  
848 and tools: a survey. *Educ Inf Technol* 22: 3121-45.
- 849 [67] Angles R, Gutierrez C (2008) Survey of graph database models. *ACM Computing*  
850 *Surveys (CSUR)* 40: 1-39.
- 851 [68] Vrandečić D, Krötzsch, M (2014) Wikidata: a free collaborative knowledgebase.  
852 *Communications of the ACM* 57: 78-85.
- 853 [69] Fox P, Hendler JA (2009) Semantic escience: encoding meaning in next-generation  
854 digitally enhanced science. Microsoft Research 2009. The Fourth Paradigm 147-  
855 152.
- 856 [70] Baker M (2016) Reproducibility crisis. *Nature* 533: 353-66.
- 857  
858  
859  
860  
861

862

863 **CRedit authorship contribution statement**

864

865 **Sergio Ruiz Saavedra:** Writing-original draft, review & editing. **Herminio García**  
866 **González:** Writing-original draft, review & editing. **Silvia Arboleya:** Writing-original  
867 draft, review & editing. **Nuria Salazar:** Writing-original draft, review & editing. **Jose**  
868 **Emilio Labra-Gayo:** Writing-original draft, review & editing. **Susana Irene Díaz:**  
869 Writing-original draft, review & editing. **Miguel Gueimonde:** Writing-original draft,  
870 review & editing. **Sonia González:** Conceptualization, Funding acquisition, Writing-  
871 original draft, review & editing. **Clara G. de los Reyes-Gavilán:** Supervision,  
872 Conceptualization, Funding acquisition, Writing-original draft, review & editing

873

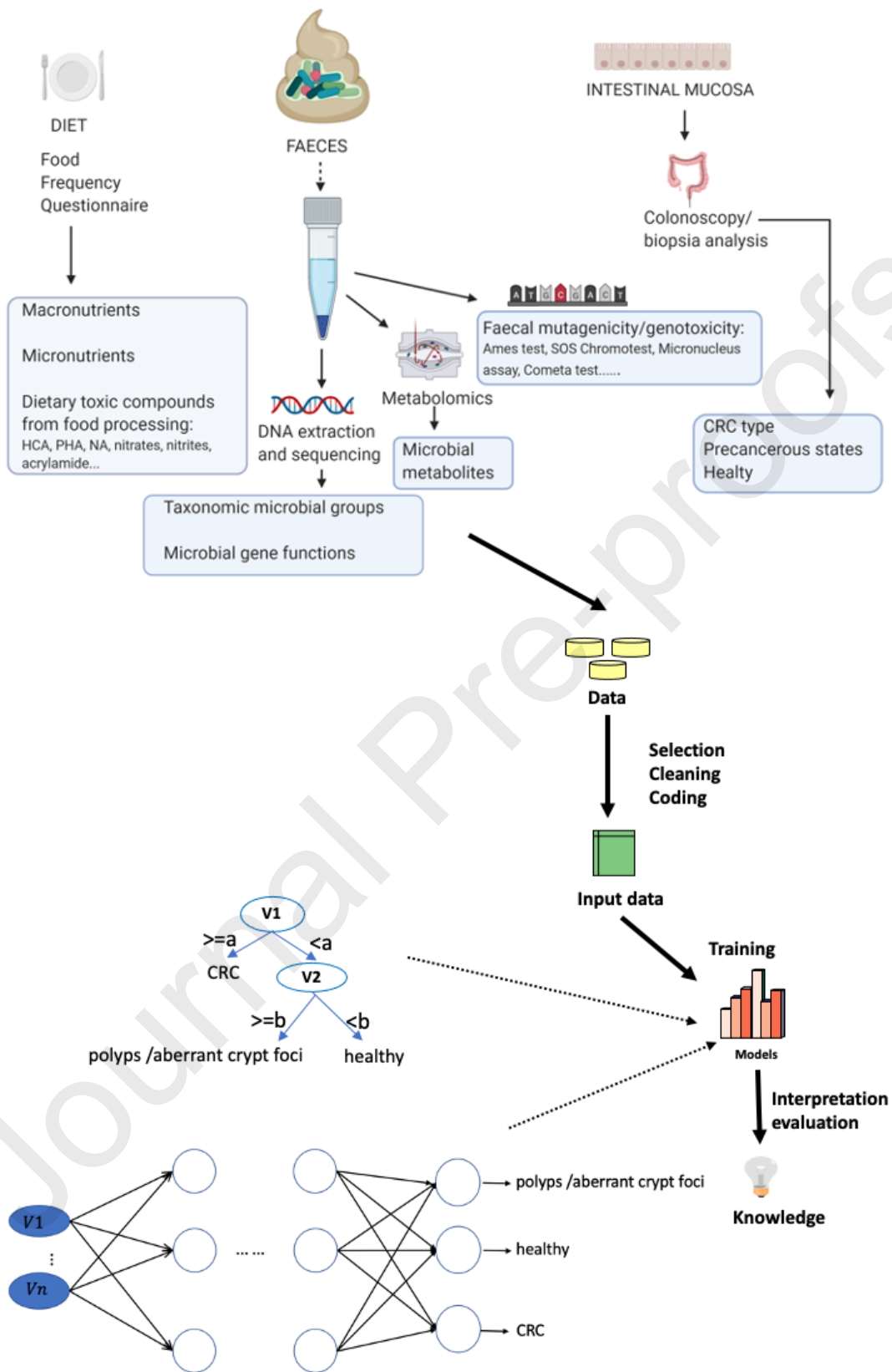
874

875 **Competing financial interest**

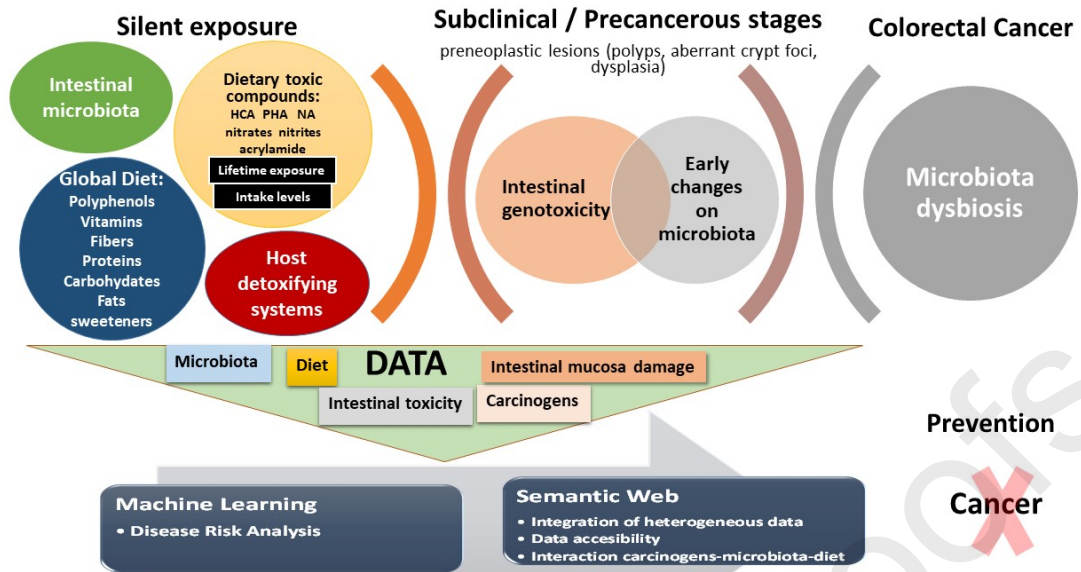
876 The authors declare no competing financial interest or personal relationships that could  
877 have influenced the content of this article.

878

879







881

Exploitation

Knowledge graph

Reasoning

To complete knowledge and to reveal not obvious present knowledge

Validation and normalisation phase

Query

To offer an engine able to solve complex queries

Isolated graphs



Data transformation phase

Publication

To enhance reusability, integration with other studies and reproducibility

Studies data



882