

1 A framework for community curation of 2 interspecies interactions literature

3

4 **Alayne Cuzick^{1,*}, James Seager¹, Valerie Wood², Martin Urban¹, Kim Rutherford² and**
5 **Kim E. Hammond-Kosack^{1,*}**

6 Author addresses:

7 ¹ Strategic area: Protecting Crops and the Environment, Rothamsted Research, Harpenden,
8 AL5 2JQ, UK

9 ² Department of Biochemistry, University of Cambridge, Cambridge, CB2 1GA, UK

10 Co-corresponding authors:

11 *To whom correspondence should be addressed. Tel: +44 1582 938240. Emails:

12 kim.hammond-kosack@rothamsted.ac.uk and alayne.cuzick@rothamsted.ac.uk

13

14 Abstract

15 The quantity and complexity of data being generated and published in biology has increased
16 substantially, but few methods exist for capturing knowledge about phenotypes derived from
17 molecular interactions between diverse groups of species, in such a way that is amenable to
18 data-driven biology and research. To improve access to this knowledge, we have
19 constructed a framework for the curation of the scientific literature studying interspecies
20 interactions, using data curated for the Pathogen-Host Interactions Database (PHI-base) as
21 a case study. The framework provides a curation tool, phenotype ontology and controlled
22 vocabularies to curate pathogen-host interaction data (at the level of the host, pathogen,
23 strain, gene and genotype). The concept of a multispecies genotype, the 'metagenotype', is
24 introduced to facilitate capturing changes in the pathogens' disease-causing abilities, and
25 host resistance or susceptibility observed by gene alterations. We report on this framework
26 and describe PHI-Canto, a community curation tool for use by publication authors.

27 Introduction

28 Recent technological advancements across the biological sciences have resulted in an
29 increasing volume of peer-reviewed publications reporting experimental data and
30 conclusions. To increase the value of this highly fragmented knowledge, biocurators
31 manually extract the data from publications and represent it in a standardized and
32 interconnected way in accordance with the FAIR (Findable, Accessible, Interoperable and
33 Reusable) Data Principles (International Society for Biocuration, 2018; Wilkinson et al.,
34 2016). The curated data is then made available in online databases, either organism- or
35 clade-specific (e.g., model organism databases) or those supporting multiple kingdoms of life
36 (e.g., PHI-base, Alliance of Genomes Resources (Alliance of Genome Resources
37 Consortium, 2020; Urban et al., 2021)). Due to the complexity of the biology, manual
38 biocuration is currently the only way to reliably represent information about function and

39 phenotype in databases and knowledge bases (Wood, Sternberg, & Lipshitz, 2022). The
40 development of curation tools with clear workflows supporting the use of biological
41 ontologies and controlled vocabularies has standardized curation efforts, reduced ambiguity
42 in annotation and improved the maintenance of the curated corpus as biological knowledge
43 evolves (International Society for Biocuration, 2018).

44 The pathogen-host interaction research communities are an example of a domain of the
45 biological sciences exhibiting a literature deluge (Figure 1). The Pathogen-Host Interactions
46 Database, PHI-base (phi-base.org), is an open access FAIR biological database containing
47 data on bacterial, fungal and protist genes proven to affect the outcome of pathogen-host
48 interactions (Rodriguez-Iglesias et al., 2016; Urban et al., 2021). Viruses are not included in
49 PHI-base. Since 2005, PHI-base has manually curated phenotype data associated with
50 underlying genome-level changes from peer-reviewed pathogen-host interaction literature.
51 Information is also provided on the target sites of some anti-infective chemistries (Urban et
52 al., 2020). This type of data is increasingly relevant, as infectious microbes continually
53 threaten global food security, human health across the life course, farmed animal health and
54 wellbeing, tree health and ecosystem resilience (Brown et al., 2012; Fisher, Hawkins,
55 Sanglard, & Gurr, 2018; Fisher et al., 2012; Smith, Machalaba, Seifman, Feferholtz, &
56 Karesh, 2019). Rising resistance to antimicrobial compounds, increased globalization, and
57 climate change indicate that infectious microbes will present ever greater economic and
58 societal threats (Bebber, Ramotowski, & Gurr, 2013; Chaloner, Gurr, & Bebber, 2021; Cook
59 et al., 2021). In order to curate relevant publications into PHI-base (version 4) professional
60 curators have, since 2011, entered 81 different data types into a text file (Urban et al., 2017).
61 However, increasing publication numbers and data complexity required more robust curation
62 procedures.

63 We were unable to locate any curation frameworks or tools capable of capturing the
64 interspecies interactions required for PHI-base. PomBase, the fission yeast
65 (*Schizosaccharomyces pombe*) database developed Canto, a web-based tool supporting

66 curation by both professional biocurators and publication authors (Rutherford, Harris, Lock,
67 Oliver, & Wood, 2014). While Canto could support annotation for multiple species, it could
68 not annotate interactions between species. Therefore, we extended and customized Canto
69 to support interspecies interactions, creating a new tool: PHI-Canto (the Pathogen-Host
70 Interaction Community Annotation Tool). Likewise, there were no existing biomedical
71 ontologies that could accurately describe pathogen-host interaction phenotypes at the depth
72 and breadth required for PHI-base. Infectious disease formation depends on a series of
73 complex and dynamic interactions between pathogenic species and their potential hosts,
74 and also requires the correct biotic and/or abiotic environmental conditions (Scholthof,
75 2007), as illustrated by the concept of the ‘disease triangle’ (Figure 2). All these interrelated
76 factors must be recorded in order to sufficiently describe a pathogen-host interaction.

77 In this study three key issues were addressed in order to develop the curation framework for
78 interspecies interactions: firstly, to support the classification of genes as ‘pathogen’ or ‘host’,
79 and enable the variations of the same gene in different strains to be captured; secondly,
80 formulating the concept of a ‘metagenotype’ to represent the interaction between specific
81 strains of both a pathogen and a host within a multispecies genotype; and thirdly, developing
82 supporting ontologies and controlled vocabularies, including the generic Pathogen-Host
83 Interaction Phenotype Ontology (PHIPO), to annotate phenotypes connected to genotypes
84 at the level of a single species (pathogen or host) and multiple species (pathogen-host
85 interaction phenotypes).

86 Results

87 Enabling multispecies curation with UniProtKB accessions

88 In any curation context, stable identifiers are required for annotated entities. The UniProt
89 Knowledgebase (UniProtKB) (UniProt Consortium, 2021) is universally recognized, provides

90 broad taxonomic protein coverage, and manually curates standard nomenclature across
91 protein families. Protein sequences are both manually and computationally annotated in
92 UniProtKB, providing a wealth of data on catalytic activities, protein structures and protein-
93 protein interactions, Gene Ontology (GO) annotations and links to PHI-base phenotypes
94 (Ashburner et al., 2000; Gene Ontology Consortium, 2021; Urban et al., 2021). To improve
95 interoperability with other resources, we used UniProtKB accession numbers for retrieving
96 protein entities, gene names and species information for display in PHI-Canto. PHI-Canto
97 accesses the UniProtKB API to automatically retrieve the entities and their associated data.

98 Developing the metagenotype to capture interspecies 99 interactions

100 To enable annotation of interspecies interactions, we developed the concept of a
101 'metagenotype', that represents the combination of a pathogen genotype and a host
102 genotype (Figure 3). A metagenotype is created after the individual genotypes from both
103 species are created. Each metagenotype can be annotated with pathogen-host interaction
104 phenotypes to capture changes in pathogenicity (caused by alterations to the pathogen) and
105 changes in virulence (caused by alterations to the host and/or the pathogen).
106 Metagenotypes must always include at least one named pathogen gene with a genotype of
107 interest, but a metagenotype can be composed from a pathogen genotype and a host
108 species (and strain) if no specific host gene is referenced in an experiment.

109 Annotation types and annotation extensions in PHI-Canto

110 In PHI-Canto, 'annotation' is the task of relating a specific piece of knowledge to a biological
111 feature. To curate a wide variety of experiment types, three groupings of annotation types
112 are available in PHI-Canto, covering 'metagenotype', 'genotype' (of a single species) and
113 'gene' annotation types (Table 1). To capture additional biologically relevant information

114 associated with an annotation, curators use annotation extensions (Huntley et al., 2014) to
115 extend the primary annotation. For the purpose of Canto and PHI-Canto, the meaning of
116 ‘annotation extension’ was broadened to capture additional properties related to the
117 annotation, such as the metagenotype used as an experimental control. The additional
118 properties that may be related to an annotation are simply referred to as ‘annotation
119 extensions’ (AEs) in this manuscript (Table 1, Supplementary file 1 and Supplementary file
120 2). Descriptions of the new AEs for PHI-Canto and the core collection of AEs from Canto are
121 available in the PHI-Canto user documentation (see the Code availability section).

122 Metagenotypes can be annotated with terms from an ontology or controlled vocabulary
123 following either the ‘pathogen-host interaction phenotype’, ‘gene-for-gene phenotype’ or
124 ‘disease name’ annotation types (Table 1). Phenotype annotations can be supported by AEs
125 providing additional qualifying information required to fully interpret the experiment, such as
126 the infected tissue of the host.

127 Phenotypes can also be curated for single species experiments, involving either the
128 pathogen or host, following the ‘single species phenotype’ annotation workflow (Table 1).
129 Single species phenotype annotations have a selection of AEs available, including the
130 protein assayed in the experiment and the severity of the observed phenotype (see example
131 from PMID:22314539 in Appendix 1).

132 PHI-Canto also supports the annotation of gene and gene product attributes to represent the
133 evolved functional role of a gene product, described here as the ‘gene annotation’ workflow
134 (Table 1). The Gene Ontology is used for annotation of a gene product’s molecular
135 functions, biological processes and cellular components, while PSI-MOD is used for the
136 annotation of protein modifications (Montecchi-Palazzi et al., 2008), and BioGRID
137 experiment types are used to capture genetic and physical interactions (Oughtred et al.,
138 2021). GO annotations are submitted to the EBI GO Annotation Database (GOA), from

139 where they are propagated to the main GO database (Gene Ontology Consortium, 2021;
140 Huntley et al., 2015).

141 Curation of interspecies interaction publications

142 Ten publications covering a wide range of typical plant, human, and animal pathogen-host
143 interactions were selected for trial curation in PHI-Canto (Table 2). These publications
144 included experiments with early acting pathogen virulence proteins, the first host targets of
145 pathogen effectors, and resistance to antifungal chemistries. These publications guided
146 development of the ontology and controlled vocabulary terms required for PHI-Canto, as well
147 as the curation methods required for different experiments. Major curation problems and
148 their solutions are summarized in Table 3, and example annotations are described below
149 and provided in Appendix 1 and Appendix 2.

150 Curating an experiment with a metagenotype

151 A large proportion of the curation in PHI-Canto requires the use of metagenotypes: one of
152 the simpler cases involves early acting virulence proteins, where a genetically modified
153 pathogen is inoculated onto a host (without a specified host gene). A metagenotype is
154 created and annotated with a phenotype term. These experiments are curated following the
155 'pathogen-host interaction phenotype' workflow, including any relevant AEs (Table 1). This
156 two-step curation process is illustrated by PMID:29020037 curation (Table 2, Appendix 1
157 and Appendix 2). The GT2 gene is deleted from the fungal plant pathogen *Zymoseptoria*
158 *septoria* and inoculated onto wheat plants; the observed phenotype 'absence of pathogen-
159 associated host lesions' (PHIPO:0000481) is annotated to the metagenotype; and the AE for
160 'infective ability' is annotated with 'loss of pathogenicity'.

161 Curating pathogen effector experiments

162 A pathogen effector is defined as an entity transferred between the pathogen and the host
163 that is known or suspected to be responsible for either activating or suppressing a host
164 process commonly involved in defense (Houterman et al., 2009; Jones & Dangl, 2006)
165 (Figure 2). To curate an effector experiment, first a metagenotype is created, then annotated
166 with a phenotype term. To indicate that the pathogen gene functions as an effector, it is
167 necessary to also make a concurrent 'gene annotation' (Table 1) with the GO biological
168 process term 'effector-mediated modulation of host process' (GO:0140418) or an
169 appropriate descendant term. This GO term has been created (with descendants) in
170 collaboration with the Gene Ontology Consortium (GOC) and is used to identify pathogen
171 effectors in PHI-base (version 5) (Supplementary file 3). Molecular functions of the pathogen
172 gene can be curated with a GO molecular function term, if reported in the literature, and
173 connected to the GO biological process term. An example of curation of a pathogen effector
174 experiment is illustrated using PMID:31804478 (Table 2 and Appendix 1) where the
175 pathogen effector Pst_12806 from *Puccinia striiformis* suppresses pattern-triggered
176 immunity in a tobacco leaf model. Here, the metagenotype is curated with the phenotype
177 'decreased level of host defense-induced callose deposition' (PHIPO:0001015) and the
178 effector is annotated with 'effector-mediated suppression of host pattern-triggered immunity'
179 (GO:0052034). A further experiment demonstrated that the pathogen effector protein was
180 able to bind to the natural host (wheat) protein PetC and inhibit the enzyme activity of PetC,
181 resulting in a GO molecular function annotation 'enzyme inhibitor activity' (GO:0004857) on
182 Pst_12806, with PetC captured as the target protein in an AE (see Appendix 1).

183 Curating experiments with a gene-for-gene relationship

184 For a gene-for-gene pathogen-host interaction type (when a known genetic interaction is
185 conferred by a specific pathogen avirulence gene product and its cognate host resistance
186 gene product) (Figure 2c, d, further described in the figure legend) (Flor, 1956; Jones &

187 Dangl, 2006; Kanyuka, Igna, Solomon, & Oliver, 2022) the ‘gene-for-gene phenotype’
188 metagenotype workflow is followed. The metagenotypes and phenotype annotations are
189 made in the same way as the standard ‘pathogen-host interaction phenotype’ workflow, but
190 with different supporting data. A new AE was developed to indicate the following three
191 components of the interaction: i) the compatibility of the interaction, ii) the functional status of
192 the pathogen gene, and iii) the functional status of the host gene. An example of an
193 annotation for a biotrophic pathogen gene-for-gene interaction has been illustrated with
194 PMID:20601497 (Table 2 and Appendix 1). Inverse gene-for-gene relationships occur with
195 necrotrophic pathogens, where the pathogen necrotrophic effector interacts with a gene
196 product from the corresponding host susceptibility locus and activates a host response that
197 benefits the pathogen (a compatible interaction). If the necrotrophic effector cannot interact
198 with the host target, then no disease occurs (an incompatible interaction) (Breen, Williams,
199 Winterberg, Kobe, & Solomon, 2016). An example of an inverse gene-for-gene interaction
200 using the appropriate AEs is illustrated with PMID:22241993 (Table 2 and Appendix 1).

201 Curating an experiment with a single species genotype in the presence
202 or absence of a chemical

203 Single species genotypes (pathogen or host) can also be annotated with phenotypes
204 following the ‘single species phenotype annotation type’ workflow (Table 1). This is
205 illustrated using PMID:22314539 in Table 2 (and Appendix 1) with an example of an *in vitro*
206 pathogen chemistry phenotype, where a single nucleotide mutation in the *Aspergillus flavus*
207 CYP51c gene confers ‘resistance to voriconazole’ (PHIPO:0000590), an antifungal agent.

208 Supporting curation of legacy information

209 PHI-Canto's curation workflows maintain support for nine high-level terms that describe
210 phenotypic outcomes essential for taxonomically diverse interspecies comparisons, which
211 were the primary annotation method used in previous versions of PHI-base (Urban et al.,

212 2015) and which are displayed in the Ensembl Genomes browser (Yates et al., 2021). For
213 example, the ‘infective ability’ AE can be used to annotate the following subset of high-level
214 terms: ‘loss of pathogenicity’, ‘unaffected pathogenicity’, ‘reduced virulence’, ‘increased
215 virulence’ and ‘loss of mutualism’ (formerly ‘enhanced antagonism’). The mapping between
216 the nine high-level terms and the PHI-Canto curation process is further described in
217 Supplementary file 3.

218 Resolving additional problems with curating complex pathogen-host 219 interactions

220 Table 3 shows a selection of the problems encountered during the development of PHI-
221 Canto and the solutions we identified. For example, recording the delivery mechanism used
222 within the pathogen-host interaction experiment. New experimental condition terms were
223 developed with a prefix of ‘delivery mechanism’, for example, ‘delivery mechanism:
224 agrobacterium’, ‘delivery mechanism: heterologous organism’, and ‘delivery mechanism:
225 pathogen inoculation’. Another issue encountered was how to record a ‘physical interaction’
226 between proteins of different species, especially for the curation of pathogen effector first
227 host targets. This was resolved by adapting the existing Canto module for curating physical
228 interactions to support two different species.

229 Development of the Pathogen-Host Interaction Phenotype 230 Ontology and additional data lists

231 To support the annotation of phenotypes in PHI-Canto, the Pathogen-Host Interaction
232 Phenotype Ontology (PHIPO) was developed. PHIPO is a species-neutral phenotype
233 ontology that describes a broad range of pathogen-host interaction phenotypes. PHIPO’s
234 terms were developed following a pre-compositional approach, where the term names and
235 semantics are composed from existing terms from other ontologies, in order to make the

236 curation process easier. For example, the curator annotates 'resistance to penicillin'
237 (PHIPO:0000692) instead of 'increased resistance to chemical' (PHIPO:0000022) and
238 'penicillin' (CHEBI:17334) separately. Terms in PHIPO have logical definitions that follow
239 design patterns from the uPheno ontology (Shefchek et al., 2020), and mapping PHIPO
240 terms to uPheno patterns is an ongoing effort. These logical definitions provide relations
241 between phenotypes in PHIPO and terms in other ontologies, such as PATO, GO, and
242 ChEBI. PHIPO is available in OWL and OBO formats from the OBO Foundry (Jackson et al.,
243 2021).

244 PHI-Canto uses additional controlled vocabularies derived from data in PHI-base. To enable
245 PHI-Canto to distinguish between pathogen and host organisms, we extracted a list of > 250
246 pathogen and > 200 host species from PHI-base (Supplementary file 4). A curated list of
247 strain names and their synonyms for the species currently curated in PHI-base was also
248 developed for use in PHI-Canto (Supplementary file 4 and 5). PHI-base uses 'strain' as a
249 grouping term for natural pathogen isolates, host cultivars and landraces, all of which are
250 included in the curated list. The curation of pathogen strain designations was motivated by
251 the NCBI Taxonomy's decision to discontinue the assignment of strain-level taxonomic
252 identifiers (Federhen et al., 2014) and a lack of standardized nomenclature for natural
253 isolates of non-model species. New strain designations can be requested by curators and
254 are reviewed by an expert prior to inclusion to ensure that each describes a novel strain
255 designation rather than a new synonym for an existing strain.

256 Annotations in PHI-Canto include experimental evidence, which is specified by a term from a
257 subset of the Evidence & Conclusion Ontology (ECO) (Giglio et al., 2019). Experimental
258 evidence codes specific to pathogen-host interaction experiments have been developed and
259 submitted to ECO. Phenotype annotations also include experimental conditions that are
260 relevant to the experiment being curated, which are sourced from the PHI-base
261 Experimental Conditions Ontology (PHI-ECO).

262 PHI-Canto includes a 'disease name' annotation type (Table 1) for annotating the name of
263 the disease caused by an interaction between the pathogen and host specified in a wild type
264 metagenotype (this annotation type is described in the PHI-Canto user documentation).
265 Diseases are specified by a controlled vocabulary of disease names (called PHIDO), which
266 was derived from disease names curated in previous versions of PHI-base.

267 Summary of the PHI-Canto curation process

268 The PHI-Canto curation process is outlined in Figure 4, Figure 4 – figure supplement 1, the
269 PHI-Canto user documentation and a detailed worked example is provided in Appendix 2.
270 Each curation session is associated with one publication (using its PubMed identifier). One
271 or more curators can collaborate on curating the same publication. An instructional email is
272 sent to curators when they begin a new curation session, and PHI-base provides further
273 guidelines on what information is needed in order to curate a publication in PHI-Canto
274 (Figure 4 – figure supplement 2) and how to identify UniProtKB accession numbers from
275 reference proteomes (Figure 4 – figure supplement 3).

276 The curator first adds genes from the publication, then creates alleles from genes,
277 genotypes from alleles and metagenotypes from pathogen and host genotypes. Pathogen
278 genotypes and host genotypes are created on separate pages, which only include genes
279 from the relevant species. A genotype can consist of multiple alleles, and a metagenotype
280 can contain multiple alleles from both the pathogen and the host. A 'copy and edit' feature
281 allows the creation of multiple similar annotations.

282 To make annotations, the curator selects a gene, genotype, or metagenotype to annotate,
283 then selects a term from a controlled vocabulary, adds experimental evidence, experimental
284 conditions, AEs (where available), and any additional comments. In PHI-Canto, the curator
285 can also specify a figure or table number from the original publication as part of the
286 annotation. Curators can use PHI-Canto's term suggestion feature to suggest new terms for

287 any controlled vocabulary in PHI-Canto, and experimental conditions can be entered as free
288 text if no suitable condition is found in PHI-ECO (new condition suggestions are reviewed
289 and approved by expert curators). The curation session can be saved and paused at various
290 stages in the entire process. Once the curation process is complete, the curator submits the
291 session for review.

292 Display and interoperability of data

293 The migration to incorporate FAIR principles fully into the PHI-base curation process will
294 promote interoperability between various data resources (Wilkinson et al., 2016). Figure 5
295 illustrates the internal and external resource dependencies for curation in PHI-Canto. URLs
296 and descriptions of the use of each resource are provided in Figure 5 – figure supplement 1.
297 All data curated in PHI-Canto will be displayed in PHI-base version 5, introduced in (Urban
298 et al., 2021). Additional detail on the data types displayed in PHI-base 5 is available in Table
299 4. Reciprocally, components of the interspecies curation framework (Figure 6a) provide data
300 to other resources (Figure 6b). For example, GO terms will be used in curation with PHI-
301 Canto and these annotations will be made available in the main GO database via the GOA
302 Database. PHI-base is a member of ELIXIR, one of the leading organizations for biological
303 resources and a major proponent of FAIR data.

304 Discussion

305 Scalable and accurate curation of data within the scientific literature is of paramount
306 importance due to the increasing quantity of publications and the complexity of experiments
307 within each publication. PHI-base is an example of a freely available, manually curated
308 database, which has been curating literature using professional curators since 2005
309 (Winnenburg et al., 2006). Here, we describe the development of PHI-Canto to allow the
310 curation of the interspecies pathogen-host interaction literature by professional curators and

311 publication authors. However, it should be noted that these developments – especially the
312 concept of annotating metagenotypes – could be of use to communities focused on different
313 types of interspecies interactions. Customizing Canto to use other ontologies and controlled
314 vocabularies is as simple as editing a configuration file, as shown in Source code 1.

315 Several adaptations to the original single species community annotation tool, Canto
316 (Rutherford et al., 2014), were required to convert this tool for interspecies use. Notably, the
317 need to annotate an interaction involving two different organisms necessitated the
318 development of a novel concept, the ‘metagenotype’, in order to record a combined
319 experimental genotype involving both a pathogen and a host. This is, to our knowledge, the
320 first example of such an approach to interspecies interaction curation.

321 Curation of pathogen-host interactions in PHI-Canto also necessitated the development of a
322 new phenotype ontology (PHIPO) to annotate pathogen-host interaction phenotypes in
323 sufficient detail across the broad range of host species that were curated in PHI-base. The
324 functional annotation of genes involved in interspecies interactions is a complex and
325 challenging task, requiring ongoing modifications to the Gene Ontology and occasionally
326 major refactoring to deprecate legacy terms (Gene Ontology Consortium, 2021). PHIPO
327 development and maintenance will also be an ongoing task, with both authors and
328 professional curators requesting new terms and edits to existing terms and the ontology
329 structure. Maintenance will be made more sustainable by the incorporation of logical
330 definitions that are aligned across phenotype ontologies in collaboration with the uPheno
331 project (Shefchek et al., 2020).

332 To improve the efficiency of the curation process, we are suggesting that authors follow an
333 author checklist during manuscript preparation (Appendix 3). This will improve the key
334 information (e.g., species names, gene identifiers etc.) in published manuscripts, thus
335 enabling more efficient comprehensive curation that is both human- and machine-readable.

336 The annotation procedures described here using PHI-Canto can be used to extract data

337 buried in small-scale publications and increase the accessibility of the curated article to a
338 wider range of potential users, for example computational biologists, thereby improving the
339 FAIR status of the data. The current data in PHI-base has been obtained from > 200 journals
340 (Figure 7) and therefore represents highly fragmented knowledge which is exceptionally
341 difficult to use by professionals in other disciplines. The feasibility of scalable community
342 curation with Canto is evidenced by PomBase, where Canto has been used by authors to
343 curate ~25% of the *S. pombe* literature, with the data being made available within 24 hours
344 of curation review and approval (<https://curation.pombase.org/pombe/stats/annotation>).

345 Future plans for PHI-Canto include addressing issues with natural sequence variation
346 between species strains. PHI-base contains information on numerous species with multiple
347 experimental strains, and natural sequence variation between strains can result in alterations
348 at the genome level that affect the subsequently observed phenotypes. Strain-specific
349 sequence variation is not captured in the reference proteomes stored by UniProt, even
350 though accession numbers from these proteomes are often used in PHI-Canto. Currently,
351 when a curator enters a gene with a taxonomic identifier below the species rank, PHI-Canto
352 maps the identifier to the corresponding identifier at the species rank (thus removing any
353 strain details from the organism name), and the curator specifies a strain to differentiate
354 gene variants in naturally occurring strains. However, this does not change the taxonomic
355 identifier linked to the UniProtKB accession number (nor its sequence), so the potential for
356 inaccuracy remains. To mitigate this, the future plan is to record the strain-specific sequence
357 of the gene using an accession number from a database from the International Nucleotide
358 Sequence Database Collaboration (Arita, Karsch-Mizrachi, & Cochrane, 2021).

359 The release of PHI-Canto to the community will occur gradually through various routes.
360 Community curation will be promoted by working with journals to capture the publication data
361 at source, at the point of manuscript acceptance. We will also target specific research
362 communities (e.g., those working on a particular pathogen) by inviting authors to curate their

363 own publications. Authors may contact us directly to request support while curating their
364 publications in PHI-Canto.

365 Methods

366 Changes to the Canto data model and configuration

367 Several new entities were added to PHI-Canto's data model in order to support pathogen-
368 host curation, as well as new configuration options (the new entities are illustrated in Figure
369 3 – figure supplement 1).

370 Pathogen and host roles

371 Genotype entities in PHI-Canto's data model were extended with an attribute indicating their
372 status as a pathogen genotype or a host genotype. Genotypes inherit their status (as
373 pathogen or host) from the organism, which in turn is classified as a pathogen or host based
374 on a configuration file that contains the NCBI Taxonomy ID (taxid) (Schoch et al., 2020) of
375 each host species in PHI-base. Only host taxids need to be specified since PHI-Canto
376 defaults to classifying a species as a pathogen if its taxid is not found in the configuration
377 file.

378 PHI-Canto also loads lists of pathogen and host species that specify the scientific name,
379 taxid, and common name (if any) of each species. These species lists are used to specify
380 which host species can be added as a component of the metagenotype in the absence of a
381 specific studied gene, and to override the scientific name provided by UniProtKB in favor of
382 the name used by the community (for example, to control whether the anamorph or
383 teleomorph name of a fungal species is displayed in PHI-Canto's user interface).

384 Metagenotype implementation

385 Metagenotypes were implemented by adding a 'metagenotype' entity to PHI-Canto's data
386 model. The metagenotype is the composition of two genotype entities. We also introduced
387 new relations into the data model to allow annotations to be related to metagenotypes
388 (previously, only genes and genotypes could be related to annotations).

389 Strain implementation

390 Support for strain curation was implemented by adding a 'strain' entity to PHI-Canto's data
391 model. Strains are related to an organism entity and its related genotype entities. In the user
392 interface, PHI-Canto uses the taxid of the organism to filter an autocomplete system, such
393 that only the strains of the specified organism are suggested. The autocomplete system can
394 also use synonyms in the strain list to suggest a strain based on its synonymous names.
395 Unknown strains are represented by a preset value of 'Unknown strain'.

396 Ontologies

397 PHIPO was developed using the Protégé ontology editor (Musen & Protégé Team, 2015).
398 PHIPO uses OBO namespaces to allow PHI-Canto to filter the terms in the ontology by
399 annotation type, ensuring that genotypes are annotated with single-species phenotypes and
400 metagenotypes with pathogen-host interaction phenotypes.

401 PHI-ECO was developed using Protégé, starting from a list of experimental conditions
402 originally developed by PomBase. PHIDO was initially derived from a list of diseases already
403 curated in PHI-base and is now maintained as a flat file that is converted into an OBO file
404 using ROBOT (Jackson et al., 2019).

405 Data availability

406 Pathogen-Host Interaction Phenotype Ontology: <http://purl.obolibrary.org/obo/phipo.owl>

- 407 PHI-base Experimental Conditions Ontology: <https://github.com/PHI-base/phi-eco>
- 408 PHIDO, the controlled vocabulary of disease names: <https://github.com/PHI-base/phido>
- 409 PHIPO Extension Ontology for gene-for-gene phenotypes: <https://github.com/PHI->
- 410 [base/phipo_ext](https://github.com/PHI-base/phipo_ext)
- 411 Location of species and strain lists used by PHI-Canto: <https://github.com/PHI-base/data>
- 412 PHI-Canto approved curation sessions (December 2022):
- 413 <https://doi.org/10.5281/zenodo.7428788>

414 Code availability

415 PHI-Canto's source code is available on GitHub, at <https://github.com/PHI-base/canto>. PHI-

416 Canto is freely licensed under the GNU General Public License version 3, with no

417 restrictions on copying, distributing, or modifying the code, for commercial use or otherwise,

418 provided any derivative works are licensed under the same terms. PHI-base provides an

419 online demo version of PHI-Canto at <https://demo-canto.phi-base.org/> which can be used for

420 evaluating the tool. The demo version and the main version of PHI-Canto will remain freely

421 available online for the foreseeable future.

422 Canto's source code is available on GitHub, at <https://github.com/pombase/canto>. Canto is

423 also freely licensed under the GNU General Public License version 3.

424 The source code for PHI-Canto's user documentation is available on GitHub, at

425 <https://github.com/PHI-base/canto-docs>. The user documentation is available online at

426 <https://canto.phi-base.org/docs/index>.

427 The source code for PHIPO is available on GitHub under a Creative Commons Attribution

428 3.0 license, at <https://github.com/PHI-base/phipo>.

429 References

- 430 Alliance of Genome Resources Consortium. (2020). Alliance of Genome Resources Portal:
431 unified model organism research platform. *Nucleic Acids Res*, *48*(D1), D650-D658.
432 doi:10.1093/nar/gkz813
- 433 Arita, M., Karsch-Mizrachi, I., & Cochrane, G. (2021). The international nucleotide sequence
434 database collaboration. *Nucleic Acids Res*, *49*(D1), D121-D124.
435 doi:10.1093/nar/gkaa967
- 436 Ashburner, M., Ball, C. A., Blake, J. A., Botstein, D., Butler, H., Cherry, J. M., . . . Sherlock,
437 G. (2000). Gene ontology: tool for the unification of biology. *Nat Genet*, *25*(1), 25-29.
438 doi:10.1038/75556
- 439 Bebber, D. P., Ramotowski, M. A. T., & Gurr, S. J. (2013). Crop pests and pathogens move
440 polewards in a warming world. *Nature Climate Change*, *3*(11), 985-988.
441 doi:10.1038/Nclimate1990
- 442 Breen, S., Williams, S. J., Winterberg, B., Kobe, B., & Solomon, P. S. (2016). Wheat PR-1
443 proteins are targeted by necrotrophic pathogen effector proteins. *Plant J*, *88*(1), 13-
444 25. doi:10.1111/tpj.13228
- 445 Brown, G. D., Denning, D. W., Gow, N. A., Levitz, S. M., Netea, M. G., & White, T. C. (2012).
446 Hidden killers: human fungal infections. *Sci Transl Med*, *4*(165), 165rv113.
447 doi:10.1126/scitranslmed.3004404
- 448 Chaloner, T. M., Gurr, S. J., & Bebber, D. P. (2021). Plant pathogen infection risk tracks
449 global crop yields under climate change. *Nature Climate Change*, *11*(8), 710+.
450 doi:10.1038/s41558-021-01104-8
- 451 Cook, N. M., Chng, S., Woodman, T. L., Warren, R., Oliver, R. P., & Saunders, D. G. (2021).
452 High frequency of fungicide resistance-associated mutations in the wheat yellow rust
453 pathogen *Puccinia striiformis f. sp. tritici*. *Pest Manag Sci*, *77*(7), 3358-3371.
454 doi:10.1002/ps.6380
- 455 Federhen, S., Clark, K., Barrett, T., Parkinson, H., Ostell, J., Kodama, Y., . . . Karsch-
456 Mizrachi, I. (2014). Toward richer metadata for microbial sequences: replacing strain-
457 level NCBI taxonomy taxids with BioProject, BioSample and Assembly records.
458 *Stand Genomic Sci*, *9*(3), 1275-1277. doi:10.4056/sigs.4851102
- 459 Fisher, M. C., Hawkins, N. J., Sanglard, D., & Gurr, S. J. (2018). Worldwide emergence of
460 resistance to antifungal drugs challenges human health and food security. *Science*,
461 *360*(6390), 739-742. doi:10.1126/science.aap7999
- 462 Fisher, M. C., Henk, D. A., Briggs, C. J., Brownstein, J. S., Madoff, L. C., McCraw, S. L., &
463 Gurr, S. J. (2012). Emerging fungal threats to animal, plant and ecosystem health.
464 *Nature*, *484*(7393), 186-194. doi:10.1038/nature10947
- 465 Flor, H. H. (1956). The complementary genic systems in Flax and Flax Rust. In M. Demerec
466 (Ed.), *Advances in Genetics* (Vol. 8, pp. 29-54): Academic Press.
- 467 Gene Ontology Consortium. (2021). The Gene Ontology resource: enriching a GOld mine.
468 *Nucleic Acids Res*, *49*(D1), D325-D334. doi:10.1093/nar/gkaa1113
- 469 Giglio, M., Tauber, R., Nadendla, S., Munro, J., Olley, D., Ball, S., . . . Chibucos, M. C.
470 (2019). ECO, the Evidence & Conclusion Ontology: community standard for evidence
471 information. *Nucleic Acids Res*, *47*(D1), D1186-D1194. doi:10.1093/nar/gky1036
- 472 Houterman, P. M., Ma, L., van Ooijen, G., de Vroomen, M. J., Cornelissen, B. J., Takken, F.
473 L., & Rep, M. (2009). The effector protein Avr2 of the xylem-colonizing fungus
474 *Fusarium oxysporum* activates the tomato resistance protein I-2 intracellularly. *Plant*
475 *J*, *58*(6), 970-978. doi:10.1111/j.1365-313X.2009.03838.x
- 476 Huntley, R. P., Harris, M. A., Alam-Faruque, Y., Blake, J. A., Carbon, S., Dietze, H., . . .
477 Mungall, C. J. (2014). A method for increasing expressivity of Gene Ontology
478 annotations using a compositional approach. *BMC Bioinformatics*, *15*, 155.
479 doi:10.1186/1471-2105-15-155

- 480 Huntley, R. P., Sawford, T., Mutowo-Meullenet, P., Shypitsyna, A., Bonilla, C., Martin, M. J.,
481 & O'Donovan, C. (2015). The GOA database: gene Ontology annotation updates for
482 2015. *Nucleic Acids Res*, 43(Database issue), D1057-1063. doi:10.1093/nar/gku1113
- 483 International Society for Biocuration. (2018). Biocuration: distilling data into knowledge.
484 *PLoS Biol*, 16(4), e2002846. doi:10.1371/journal.pbio.2002846
- 485 Jackson, R., Matentzoglou, N., Overton, J. A., Vita, R., Balhoff, J. P., Buttigieg, P. L., . . .
486 Peters, B. (2021). OBO Foundry in 2021: operationalizing open data principles to
487 evaluate ontologies. *Database (Oxford)*, 2021. doi:10.1093/database/baab069
- 488 Jackson, R. C., Balhoff, J. P., Douglass, E., Harris, N. L., Mungall, C. J., & Overton, J. A.
489 (2019). ROBOT: a tool for automating ontology workflows. *BMC Bioinformatics*,
490 20(1), 407. doi:10.1186/s12859-019-3002-3
- 491 Jones, J. D., & Dangl, J. L. (2006). The plant immune system. *Nature*, 444(7117), 323-329.
492 doi:10.1038/nature05286
- 493 Kanyuka, K., Igna, A. A., Solomon, P. S., & Oliver, R. P. (2022). The rise of necrotrophic
494 effectors. *New Phytol*, 233(1), 11-14. doi:10.1111/nph.17811
- 495 Montecchi-Palazzi, L., Beavis, R., Binz, P. A., Chalkley, R. J., Cottrell, J., Creasy, D., . . .
496 Garavelli, J. S. (2008). The PSI-MOD community standard for representation of
497 protein modification data. *Nat Biotechnol*, 26(8), 864-866. doi:10.1038/nbt0808-864
- 498 Musen, M. A., & Protege Team. (2015). The Protege project: A look back and a look forward.
499 *AI Matters*, 1(4), 4-12. doi:10.1145/2757001.2757003
- 500 Oughtred, R., Rust, J., Chang, C., Breitkreutz, B. J., Stark, C., Willems, A., . . . Tyers, M.
501 (2021). The BioGRID database: A comprehensive biomedical resource of curated
502 protein, genetic, and chemical interactions. *Protein Sci*, 30(1), 187-200.
503 doi:10.1002/pro.3978
- 504 Rodriguez-Iglesias, A., Rodriguez-Gonzalez, A., Irvine, A. G., Sesma, A., Urban, M.,
505 Hammond-Kosack, K. E., & Wilkinson, M. D. (2016). Publishing FAIR data: an
506 exemplar methodology utilizing PHI-base. *Front Plant Sci*, 7, 641.
507 doi:10.3389/fpls.2016.00641
- 508 Rutherford, K. M., Harris, M. A., Lock, A., Oliver, S. G., & Wood, V. (2014). Canto: an online
509 tool for community literature curation. *Bioinformatics*, 30(12), 1791-1792.
510 doi:10.1093/bioinformatics/btu103
- 511 Schoch, C. L., Ciufo, S., Domrachev, M., Hottton, C. L., Kannan, S., Khovanskaya, R., . . .
512 Karsch-Mizrachi, I. (2020). NCBI Taxonomy: a comprehensive update on curation,
513 resources and tools. *Database (Oxford)*, 2020. doi:10.1093/database/baaa062
- 514 Scholthof, K. B. (2007). The disease triangle: pathogens, the environment and society. *Nat*
515 *Rev Microbiol*, 5(2), 152-156. doi:10.1038/nrmicro1596
- 516 Shefchek, K. A., Harris, N. L., Gargano, M., Matentzoglou, N., Unni, D., Brush, M., . . . Osumi-
517 Sutherland, D. (2020). The Monarch Initiative in 2019: an integrative data and
518 analytic platform connecting phenotypes to genotypes across species. *Nucleic Acids*
519 *Res*, 48(D1), D704-D715. doi:10.1093/nar/gkz997
- 520 Smith, K. M., Machalaba, C. C., Seifman, R., Feferholtz, Y., & Karesh, W. B. (2019).
521 Infectious disease and economics: The case for considering multi-sectoral impacts.
522 *One Health*, 7, 100080. doi:10.1016/j.onehlt.2018.100080
- 523 UniProt Consortium. (2021). UniProt: the universal protein knowledgebase in 2021. *Nucleic*
524 *Acids Res*, 49(D1), D480-D489. doi:10.1093/nar/gkaa1100
- 525 Urban, M., Cuzick, A., Rutherford, K., Irvine, A., Pedro, H., Pant, R., . . . Hammond-Kosack,
526 K. E. (2017). PHI-base: a new interface and further additions for the multi-species
527 pathogen-host interactions database. *Nucleic Acids Res*, 45(D1), D604-D610.
528 doi:10.1093/nar/gkw1089
- 529 Urban, M., Cuzick, A., Seager, J., Wood, V., Rutherford, K., Venkatesh, S. Y., . . .
530 Hammond-Kosack, K. E. (2020). PHI-base: the pathogen-host interactions database.
531 *Nucleic Acids Res*, 48(D1), D613-D620. doi:10.1093/nar/gkz904
- 532 Urban, M., Cuzick, A., Seager, J., Wood, V., Rutherford, K., Venkatesh, S. Y., . . .
533 Hammond-Kosack, K. E. (2021). PHI-base in 2022: a multi-species phenotype

534 database for Pathogen-Host Interactions. *Nucleic Acids Res.*
535 doi:10.1093/nar/gkab1037
536 Urban, M., Pant, R., Raghunath, A., Irvine, A. G., Pedro, H., & Hammond-Kosack, K. E.
537 (2015). The Pathogen-Host Interactions database (PHI-base): additions and future
538 developments. *Nucleic Acids Res*, 43(Database issue), D645-655.
539 doi:10.1093/nar/gku1165
540 Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., . . .
541 Mons, B. (2016). The FAIR Guiding Principles for scientific data management and
542 stewardship. *Sci Data*, 3, 160018. doi:10.1038/sdata.2016.18
543 Winnenburger, R., Baldwin, T. K., Urban, M., Rawlings, C., Kohler, J., & Hammond-Kosack, K.
544 E. (2006). PHI-base: a new database for pathogen host interactions. *Nucleic Acids*
545 *Res*, 34(Database issue), D459-464. doi:10.1093/nar/gkj047
546 Wood, V., Sternberg, P. W., & Lipshitz, H. D. (2022). Making biological knowledge useful for
547 humans and machines. *Genetics*, 220(4). doi:10.1093/genetics/iyac001
548 Yates, A. D., Allen, J., Amode, R. M., Azov, A. G., Barba, M., Becerra, A., . . . Flicek, P.
549 (2021). Ensembl Genomes 2022: an expanding genome resource for non-
550 vertebrates. *Nucleic Acids Res.* doi:10.1093/nar/gkab1007

551

552

553 Acknowledgements

554 We thank former post-doctoral PHI-base team member Dr Alistair Irvine for adding chemical
555 entries to ChEBI. Dr Paul Kersey, formerly the non-vertebrate Ensembl team leader, is
556 thanked for helpful discussions and ideas on community engagement. We thank Dr Midori
557 Harris (formerly of University of Cambridge, UK) for providing valuable input into the
558 development of PHIPO based on her extensive knowledge of FYPO. Dr Pascale Gaudet
559 (Swiss-Prot, Swiss Institute of Bioinformatics) is thanked for the generation and editing of
560 GO terms involved in interspecies interactions. We also thank Drs Chris Stephens and Ana
561 Machado-Wood (both formerly of Rothamsted Research) for completing the trial curation of
562 articles into beta versions of PHI-Canto and providing invaluable feedback and suggestions
563 for further improvements. The Molecular Connections team based in Bangalore India while
564 developing the PHI-base 5 website, provided useful feedback on data interoperability
565 between PHI-Canto and the new gene-centric version of PHI-base.

566 Funding

567 PHI-base is funded by the UK Biotechnology and Biological Sciences Research Council
568 (BBSRC) Grants BB/S020020/1 and BB/S020098/1. Rothamsted authors M.U., and K.H.K.
569 receive additional BBSRC grant-aided support as part of the Institute Strategic Programme
570 Designing Future Wheat Grant (BB/P016855/1). This work was conducted using the Protégé
571 resource, which is supported by grant GM10331601 from the National Institute of General
572 Medical Sciences of the United States National Institutes of Health.

573 Author's contributions

574 AC wrote the initial manuscript draft. JS, VW, KR, MU and KHK provided comments on
575 various manuscript versions. AC, JS, MU and KHK prepared the figures and tables. AC and
576 JS prepared the supplementary files.

577 Corresponding authors

578 Correspondence to kim.hammond-kosack@rothamsted.ac.uk or
579 alayne.cuzick@rothamsted.ac.uk.

580 Ethics declarations

581 Competing interests

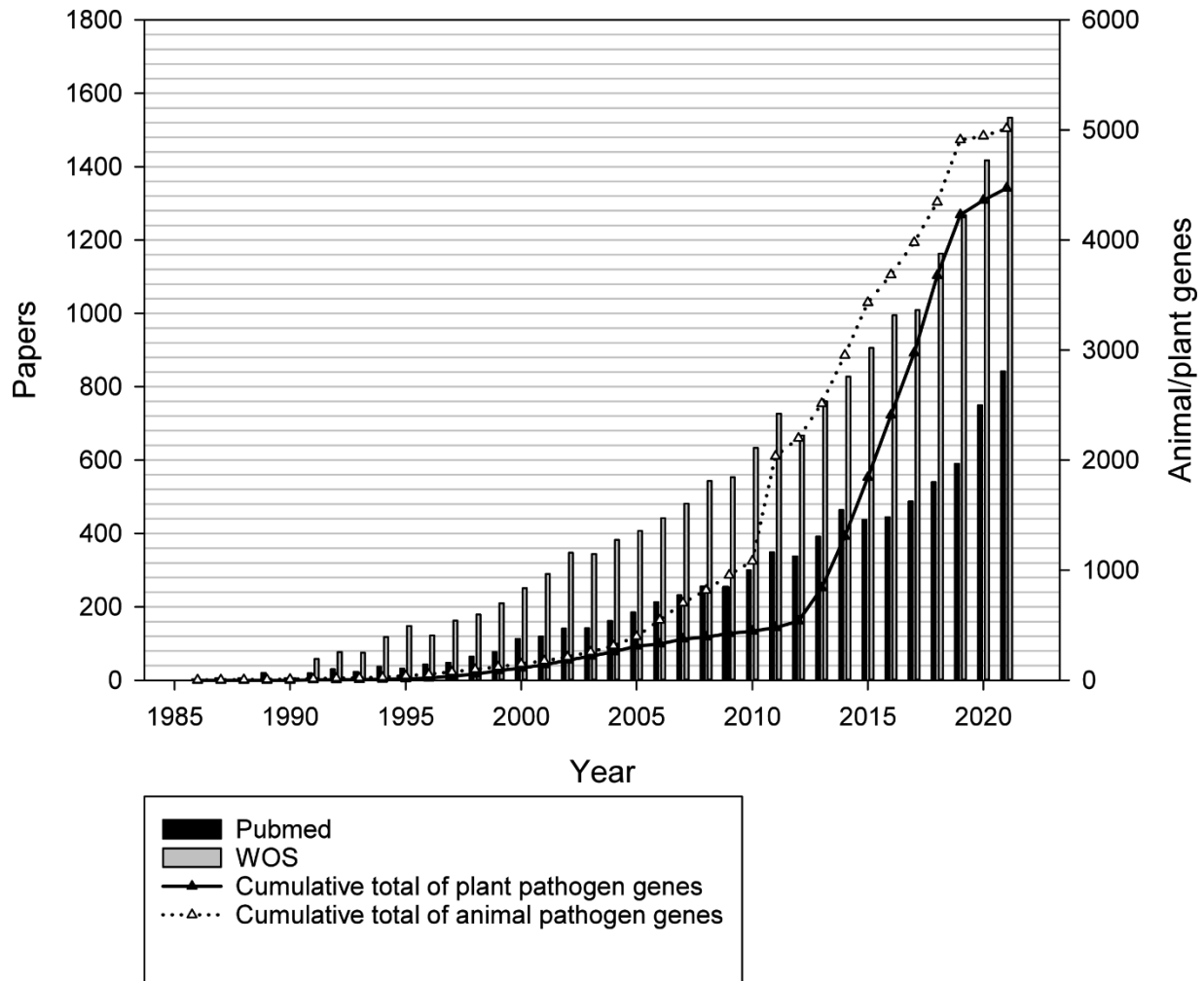
582 The authors declare no competing interests.

583

584 Tables and Figures

585

586



587

588

589 **Figure 1. Increase of molecular host-pathogen interaction publications and gene phenotype**

590 **information during the last 35 years curated in PHI-base.** Grey bars show the number of

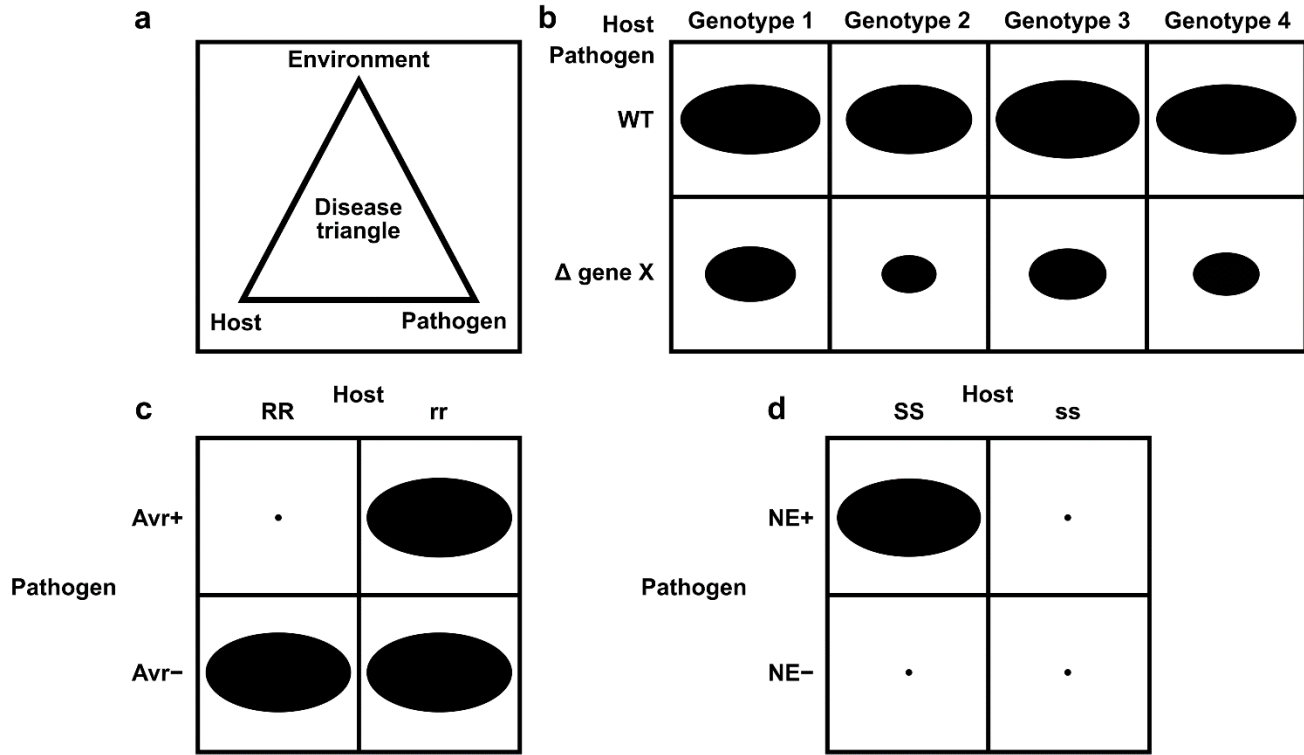
591 publications in the Web of Science Core Collection database retrieved with search term "(fung* or

592 yeast) and (gene or factor) and (pathogenicity or virulen* or avirulence gene*)". Black vertical bars

593 show the number of articles retrieved from PubMed (searching on title and abstract). Black and white

594 triangles show the number of curated animal and plant pathogen genes, respectively.

595



596
597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

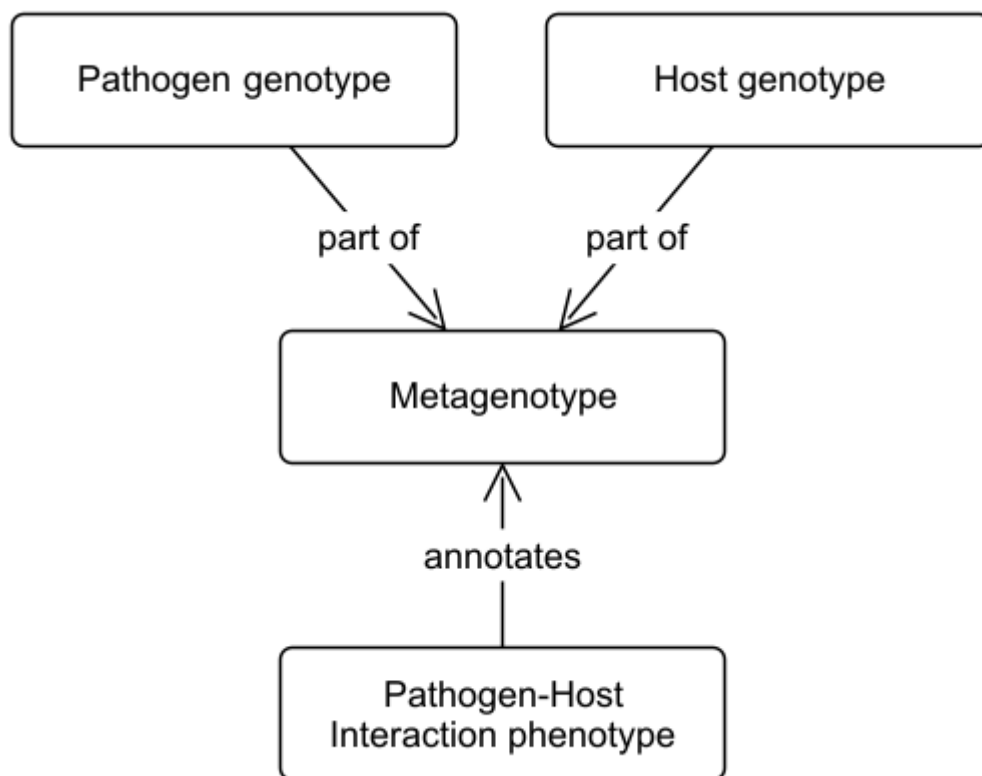
614

615

616

Figure 2. Schematic representation of pathogen-host interactions. (a) the disease triangle illustrates the requirement for the correct abiotic and biotic environmental conditions to ensure disease when an adapted pathogen encounters a suitable host; (b) a non gene-for-gene genetic relationship where compatible interactions result in disease on all host genotypes (depicted as genotypes 1–4), but the extent of disease formation is influenced to a greater or lesser extent by the presence or absence of a single pathogen virulence gene product X. In host genotypes 1 and 3, the pathogen gene product X is the least required for disease formation. The size of each black oval in each of the eight genetic interactions indicates the severity of the disease phenotype observed, with a larger oval indicating greater severity; (c) a gene-for-gene genetic relationship. In this genetic system, considerable specificity is observed, which is based on the direct or indirect interaction of a pathogen avirulence (*Avr*) effector gene product with a host resistance (*R*) gene product to determine specific recognition (an incompatible interaction), which is typically observed in biotrophic interactions ((Jones & Dangl, 2006)). In one scenario, the product of the *Avr* effector gene binds to the product of the *R* gene (a receptor) to activate host resistance mechanisms. In another scenario, the product of the *Avr* effector gene binds to an essential host target which is guarded by the product of the *R* gene (a receptor). Once *Avr* effector binding is detected, host resistance mechanisms are activated. The absence of the *Avr* effector product or the absence of the *R* gene product leads to susceptibility (a compatible interaction). The small black dot indicates no disease formation, and the large black oval indicates full disease formation, and (d) an inverse gene-for-gene genetic relationship. Again, considerable specificity is observed based on the interaction of a pathogen necrotrophic effector (*NE*) with a host susceptibility (*S*) target to determine specific recognition. The product of the pathogen *NE* gene binds to the product of the *S* gene (a receptor) to activate host susceptibility mechanisms.

617



618

619

620

621

622

623

Figure 3. Conceptual model showing the relationship between metagenotypes, genotypes and annotations. The curator selects a pathogen genotype and a host genotype to combine into a metagenotype. The metagenotype can be annotated with pathogen-host interaction phenotypes from PHIPO (the Pathogen-Host Interaction Phenotype Ontology).

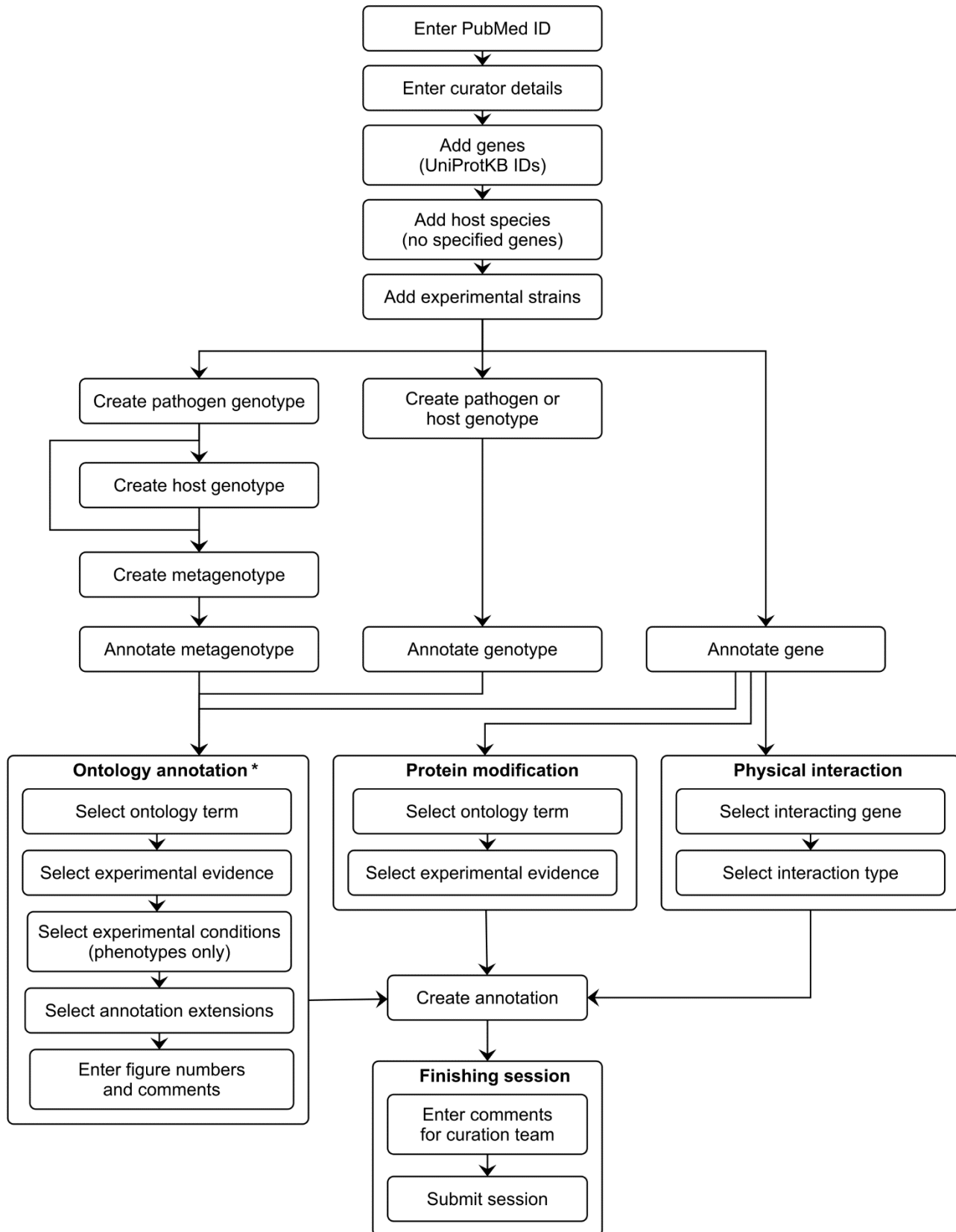
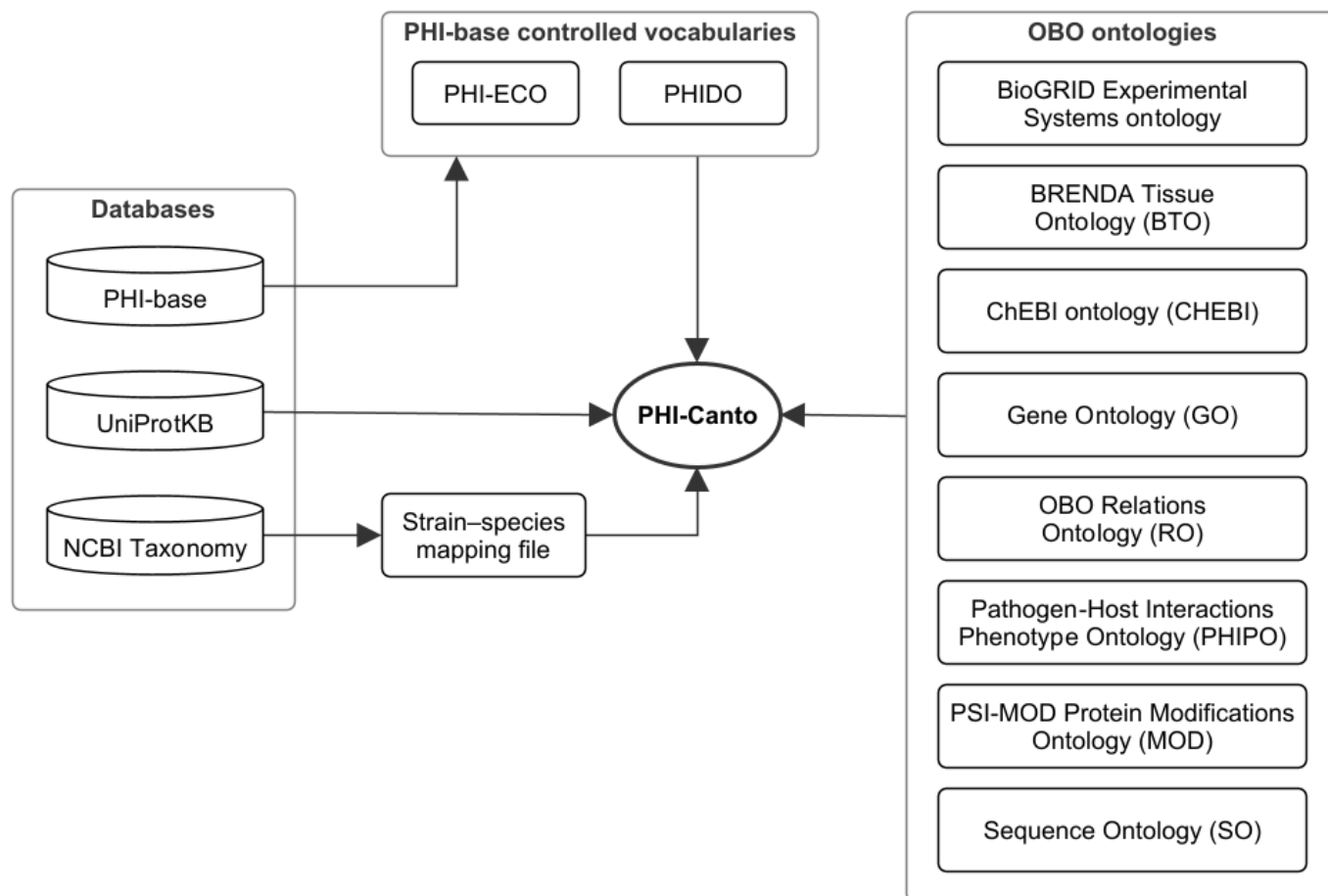


Figure 4. PHI-Canto curation workflow diagram. This diagram shows the curation workflow from the start of a curation session to its submission. The PubMed ID of the publication to be curated is entered and the title is automatically retrieved. The curator enters their name, email address and ORCID iD. On the species and genes page, the experimental pathogen and host genes are entered using UniProtKB accession numbers, and for experiments where a mutant pathogen genotype is assayed on a wild type host with no specified genes, there is the option to select the host species from an autocomplete menu. Information on the specific experimental strains used for each species is entered. After entering this initial information, the curator follows one of three distinct workflows depending on the biological feature the user wants to annotate (metagenotype, genotype or gene annotation type). Except for genes, biological features are created by composing less complex features: genotypes from alleles (generated in the pathogen or host genotype management pages), and metagenotypes from genotypes (generated in the metagenotype management page). Biological features are annotated with terms from a controlled vocabulary (usually an ontology),

636 plus additional information that varies based on the annotation type. The curator has the option to generate further
637 annotations after creating one, but this iterative process is not represented in the diagram for the sake of brevity. After
638 all annotations have been made, the session is submitted to PHI-base. * Note that the 'Ontology annotation' group
639 covers multiple annotation types, all of which annotate biological features with terms from an ontology or controlled
640 vocabulary. These annotation types are described in Table 1.
641



642
643 **Figure 5. Network diagram showing the data resources used by PHI-Canto.** Of the databases shown, PHI-base
644 provides data (experimental conditions and disease names) used to create terms in the PHI-base controlled
645 vocabularies; UniProtKB provides accession numbers for proteins that PHI-Canto uses to identify genes; and the
646 NCBI Taxonomy database is used to generate a mapping file relating taxonomic identifiers lower than species rank to
647 their nearest taxonomic identifiers at species rank. The OBO ontologies group contains ontologies in the OBO format
648 that PHI-Canto uses for its annotation types. The parenthesized text after the ontology name indicates the term prefix
649 for the ontology.
650

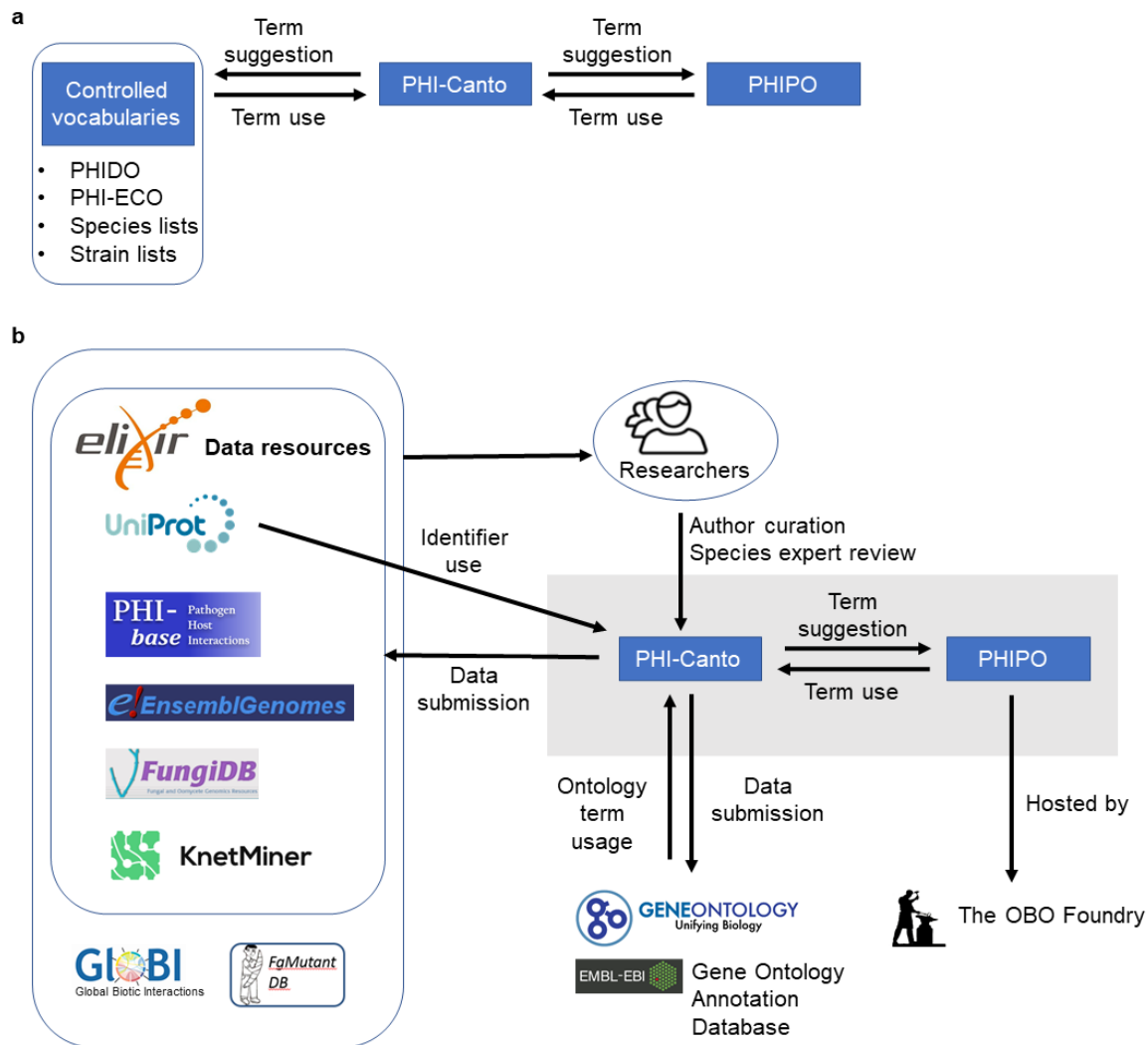
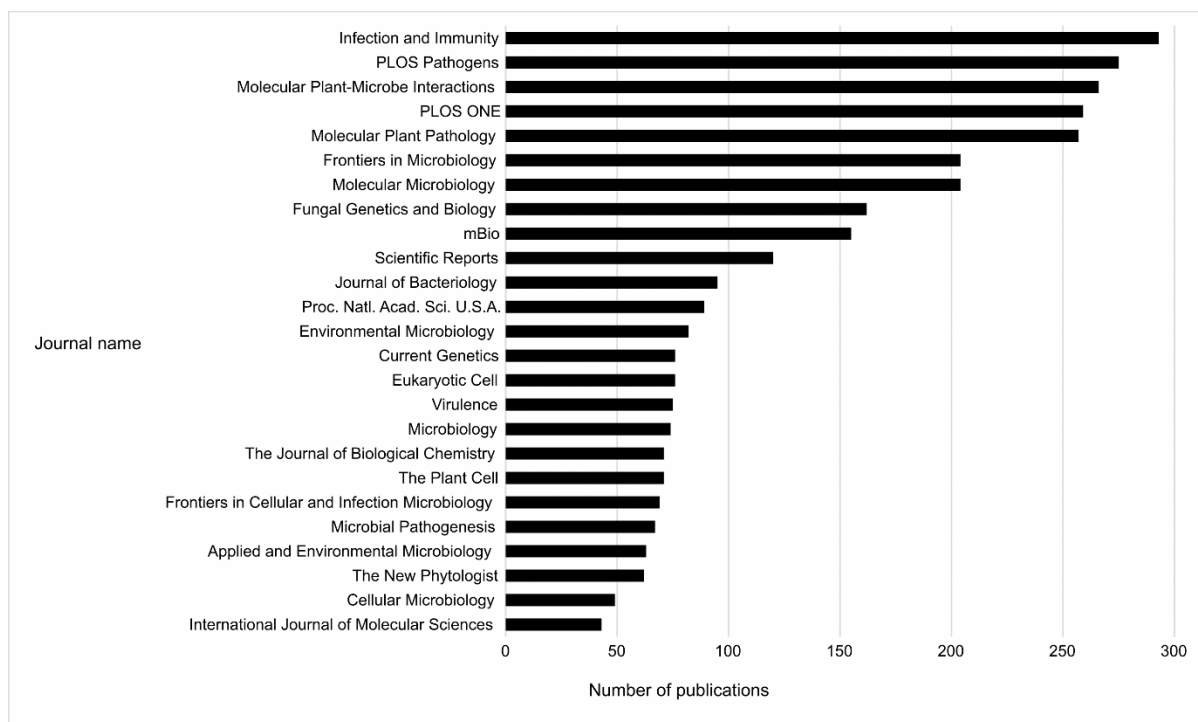


Figure 6. The interspecies curation framework and the interoperability of PHI-Canto.

(a) The interspecies curation framework consists of three main components. Firstly, a curation tool called PHI-Canto (The Pathogen Host Interaction Community Annotation Tool), secondly, a new species neutral phenotype ontology called PHIPO (the Pathogen-Host Interaction Phenotype Ontology), and thirdly, a selection of additional controlled vocabularies for disease names (PHIDO), experimental conditions (PHI-ECO), pathogen and host species, and natural strains associated with each species. The two-way arrows indicate that terms from the ontology and controlled vocabularies are used in curation with PHI-Canto, and that new terms required for curation may be suggested for inclusion within the ontology and controlled vocabularies. (b) The PHI-Canto and PHIPO content curation framework (grey box) uses persistent identifiers and cross-referenced information from UniProt, Ensembl Genomes and the Gene Ontology. PHIPO is made available at the OBO Foundry. Newly minted wild type gene annotations are suggested for inclusion into the Gene Ontology via the EBI Gene Ontology Annotation database. Data curated in PHI-Canto will be shared with ELIXIR data resources such as UniProtKB, Ensembl Genomes, FungiDB, and KnetMiner, and will be provided on request to other databases (FgMutantDB, GloBI). Researchers can look up curated information via the PHI-base web interface or can download the whole dataset for inclusion in their bioinformatics pipelines. Authors can submit data to PHI-base by curating their publications into PHI-Canto. The origin of data is indicated by directional arrows.



669

670

671

672

673

674

Figure 7. Top 25 Journals in PHI-base.

Bar chart showing the top 25 journals by number of publications curated in PHI-base, as of version 4.13 (published 9 May 2022). Publication counts were generated by extracting every unique PubMed identifier (PMID) from PHI-base, then using the Entrez Programming Utilities (E-Utilities) to retrieve the journal name for each PMID, and finally summing the count of journal names.

675
676
677

Table 1. Annotation types and selected annotation extensions used in PHI-Canto.

Annotation type	Annotation extensions ¹	Annotation value
Gene annotation types ²		
Gene Ontology annotation		Gene Ontology term
	with host species	NCBI Taxonomy ID
	with symbiont species	NCBI Taxonomy ID
Wild type expression		PomBase Gene Expression ontology term
	during	Gene Ontology biological process term ³
	in presence of	Chemical entity (ChEBI ontology)
	tissue type	BRENDA Tissue Ontology term
Genotype annotation types		
Single species phenotype (Pathogen phenotype and Host phenotype)		PHIPO term (single-species phenotype branch)
	affected proteins	UniProtKB accession number (one for each affected protein)
	assayed RNA	UniProtKB accession number
	assayed protein	UniProtKB accession number
	observed in organ	BRENDA Tissue Ontology term ⁴
	penetrance	Qualitative value (low, normal, high, complete) or quantitative value (percentage)
	severity	Qualitative value (low, normal, high, variable) or quantitative value (percentage)
Metagenotype annotation types		
Pathogen-host interaction phenotype or Gene-for-gene phenotype		PHIPO term (pathogen-host interaction phenotype branch)
	affected proteins	UniProtKB accession number (one for each affected protein)
	assayed protein	UniProtKB accession number
	assayed RNA	UniProtKB accession number
	compared to control metagenotype	Metagenotype ⁵
	extent of infectivity ⁶	PHIPO term
	gene-for-gene interaction ⁷	PHIPO Extension (PHIPO_EXT) ontology term
	host tissue infected	BRENDA Tissue Ontology term
	inverse gene-for-gene interaction ⁷	PHIPO Extension (PHIPO_EXT) ontology term
	outcome of interaction ⁶	PHIPO term
	penetrance	Qualitative value (low, normal, high, complete) or quantitative value (percentage)
	severity	Qualitative value (low, normal, high) or quantitative value (percentage)
Disease name		PHIDO term ⁸
	host tissue infected	BRENDA Tissue Ontology term

678
679
680
681
682
683
684
685
686

¹ PHI-Canto uses 44 annotation extension (AE) relations, of which 9 are unique to PHI-base, while the remaining 35 are shared with PomBase.

² Additional AEs shared with PomBase for the gene annotation types are available in Supplementary file 2.

³ Restricted to GO:0022403, GO:0033554, GO:0072690, GO:0051707 and their descendant terms.

⁴ Restricted to BTO:0001489, BTO:0001494, BTO:0001461 and their descendant terms.

⁵ Metagenotypes are selected from those already added to the curation session.

⁶ AE only applies to pathogen-host interaction phenotypes.

⁷ AE only applies to gene-for-gene phenotypes. ⁸ Curated list of disease names.

687
688

Table 2. Publications selected for trial curation using PHI-Canto.

Subject of publication	PMID	Publication title	Genotype ¹⁰ annotated with	Metagenotype ¹¹ annotated with
Bacteria-human interaction	28715477 ¹	The RhlR quorum-sensing receptor controls <i>Pseudomonas aeruginosa</i> pathogenesis and biofilm development independently of its canonical homoserine lactone autoinducer.	Pathogen phenotype	unaffected pathogenicity altered pathogenicity or virulence
Fungal-human interaction/novel antifungal target	28720735 ²	A nonredundant phosphopantetheinyl transferase, PptA, is a novel antifungal target that directs secondary metabolite, siderophore, and lysine biosynthesis in <i>Aspergillus fumigatus</i> and is critical for pathogenicity.	Pathogen phenotype	unaffected pathogenicity altered pathogenicity or virulence
Secondary metabolite clusters required for pathogen virulence	30459352 ²	Phosphopantetheinyl transferase (Ppt)-mediated biosynthesis of lysine, but not siderophores or DHN melanin, is required for virulence of <i>Zymoseptoria tritici</i> on wheat.	Pathogen phenotype	unaffected pathogenicity altered pathogenicity or virulence
Early acting virulence proteins	29020037 ^{2,3}	A conserved fungal glycosyltransferase facilitates pathogenesis of plants by enabling hyphal growth on solid surfaces.	Pathogen phenotype	altered pathogenicity or virulence
Mutualism interaction	16517760 ⁴	Reactive oxygen species play a role in regulating a fungus-perennial ryegrass mutualistic interaction	Pathogen phenotype	mutualism
First host targets of pathogen effectors	31804478 ^{2,5}	An effector protein of the wheat stripe rust fungus targets chloroplasts and suppresses chloroplast function.	N/A	altered pathogenicity or virulence a pathogen effector
Receptor decoys	30220500 ⁵	Suppression of plant immunity by fungal chitinase-like effectors.	Pathogen phenotype	a pathogen effector
R-Avr interactions	20601497 ^{6,7}	Activation of an Arabidopsis resistance protein is specified by the <i>in planta</i> association of its leucine-rich repeat domain with the cognate oomycete effector.	Host phenotype	a pathogen effector a gene-for-gene interaction
Fungal toxins required for virulence on plants	22241993 ⁸	The cysteine rich necrotrophic effector SnTox1 produced by <i>Stagonospora nodorum</i> triggers susceptibility of wheat lines harboring Snn1.	N/A	a pathogen effector a gene-for-gene interaction (inverse)
Resistance to antifungal chemistries	22314539 ⁹	The T788G mutation in the cyp51C gene confers voriconazole resistance in <i>Aspergillus flavus</i> causing aspergillosis.	Pathogen phenotype Pathogen chemistry phenotype	

689
690
691
692
693
694
695
696
697
698
699
700
701
702

¹ Example of curating 'unaffected pathogenicity' available in Appendix 1.

² Example of curating 'altered pathogenicity or virulence' available in Appendix 1 and Appendix 2.

³ Example of '*in vitro* pathogen phenotype' available in Appendix 1.

⁴ Example of curating 'mutualism' available in Appendix 1.

⁵ Example of curating 'a pathogen effector' available in Appendix 1.

⁶ Example of curating 'a gene-for-gene interaction' available in Appendix 1.

⁷ Example of '*in vitro* host phenotype' available in Appendix 1.

⁸ Example of curating 'an inverse gene-for-gene interaction' available in Appendix 1.

⁹ Example of '*in vitro* pathogen chemistry phenotype' available in Appendix 1.

¹⁰ Single species genotypes could be annotated with either a pathogen phenotype, a pathogen chemistry phenotype, or a host phenotype. Genotypes are annotated with *in vitro* or *in vivo* phenotypes from PHIPO, using either the Pathogen phenotype or Host phenotype annotation type workflow.

¹¹ Metagenotype comprises of a pathogen and a host genotype in combination. Phenotypes from PHIPO can be annotated to metagenotypes using either the 'Pathogen-Host Interaction Phenotype' or 'Gene-for-Gene Phenotype' annotation type workflow.

703
704

Table 3. Issues encountered whilst curating ten example publications with PHI-Canto.

Curated feature	Problem description	Solution	Context in PHI-Canto	Example
Species strain	UniProtKB sequence information is commonly from a reference genome strain. This sequence may differ from the experimental strain curated in PHI-Canto.	Develop a selectable list of strains for curators to assign to the genotype (and metagenotype).	Strain selected after UniProtKB entry on gene entry page. Strain used within genotype creation.	URL ¹ All phenotype annotation examples in Appendix 1 contain a 'strain name' within the genotype / metagenotype.
Delivery mechanism	Pathogen-host interaction experiments use a wide array of mechanisms to deliver the treatment of choice (to cells, tissues, and host and non-host species) which are required for experimental interpretation.	Develop terms prefixed with 'delivery mechanism' in the Pathogen-Host Interaction Experimental Conditions Ontology (PHI-ECO).	Selection of experimental conditions whilst making a phenotype annotation to a metagenotype.	URL ² Examples in Appendix 1 PMID:20601497, PMID:31804478 and PMID:22241993.
Physical interaction	Physical interactions (i.e., protein-protein interactions) could only be annotated between proteins of the same species, so it was not possible to annotate interactions between a pathogen effector and its first host target.	Adapt the 'Physical Interaction' annotation type to store gene and species information from two organisms (instead of one).	Physical Interaction annotation type.	URL ³
Pathogen effector	There was no available ontology term to describe a 'class' pathogen effector (a 'transferred entity from pathogen to host'), because effectors have heterogeneous functions (specific enzyme inhibitors, modulating host immune responses, and targeting host gene-silencing mechanisms). Effector is not a phenotype, and so did not fit into the Pathogen-Host Interaction Phenotype Ontology (PHIPO).	Develop new Gene Ontology (GO) biological process terms (and children), to group 'effector-mediated' processes.	GO Biological Process annotation on a pathogen gene.	URL ⁴ Example in Appendix 1 PMID:31804478.
Wild type control phenotypes	Natural sequence variation between strains of both pathogen and host organisms can alter the phenotypic outcome within an interaction. The wild type metagenotype phenotype needs to be curated so that the phenotype of an altered metagenotype is informative.	Allow creation of metagenotypes containing wild type genes. Develop a new annotation extension (AE) property 'compared to control', used in annotation of altered metagenotypes.	Annotation of phenotypes and AEs to metagenotypes (using the 'PHI phenotype' or 'Gene for Gene phenotype' annotation type).	URL ⁵ Examples in Appendix 1 PMID:28715477, PMID:16517760, PMID:29020037, PMID:20601497, PMID:22241993.
Chemistry	How to record chemicals for resistance or sensitivity phenotypes.	Follow PomBase model to pre-compose PHIPO terms to include chemical names from the ChEBI ontology.	Annotation of phenotypes to single species genotypes.	URL ⁴ Example in Appendix 1 PMID:22314539.
Gene for gene interactions	Complex gene-for-gene interactions within plant pathogen-host interactions required additional detail to describe the function of the pathogen and host genes within the metagenotype (including the specified strains).	Develop the additional metagenotype curation type 'Gene for Gene Phenotype'. Develop two new AEs, 'gene_for_gene interaction' and 'inverse gene_for_gene interaction', using PHIPO_EXT terms describing three components of the interaction.	Annotation of phenotypes and AEs to metagenotypes using the 'Gene for Gene Phenotype' annotation type.	URL ⁴ Examples in Appendix 1 PMID:20601497 and PMID:22241993.
Nine high-level legacy terms (from PHI-base 4)	PHI-base should incorporate legacy data from PHI-base 4 into new PHI-base 5 gene-centric pages.	Maintain the nine high level terms as 'tags' within the new PHI-base 5 user interface. Develop mapping methods to enable this.	Three locations described in Supplementary file 3.	Urban et al., 2015 NAR (PMID:25414340).

705
706
707
708
709

*Namely, i) the compatibility of the interaction ii) the functional status of the pathogen gene and iii) the functional status of the host gene.
 URL¹ https://canto.phi-base.org/docs/getting_started#adding_strains, URL² https://canto.phi-base.org/docs/phipo_annotation#experimental_conditions, URL³ https://canto.phi-base.org/docs/physical_interaction_annotation,
 URL⁴ https://canto.phi-base.org/docs/phipo_annotation#pathogen_host_interaction_phenotypes, URL⁵ https://canto.phi-base.org/docs/genotypes#metagenotype_management

710 **Table 4.** Automatically and manually curated types of data displayed in gene-centric PHI-base 5.
711

Data type	Data source
Metadata	
Entry Summary ¹	UniProtKB ²
Pathogen species	NCBI Taxonomy ²
Pathogen strain	PHI-base strain list
Host species	NCBI Taxonomy ²
Host strain	PHI-base strain list
Publication	PubMed ²
Phenotype annotation sections	
Pathogen-Host Interaction Phenotype	PHIPO ³ pathogen-host interaction phenotype branch
Gene-for-Gene Phenotype	PHIPO pathogen-host interaction phenotype branch
Pathogen Phenotype	PHIPO single species phenotype branch
Host Phenotype	PHIPO single species phenotype branch
Other annotation sections	
Disease name	PHIDO
GO Molecular Function	GO ⁴
GO Biological Process	GO
GO Cellular Component	GO
Wild type RNA level	FYPO_EXT ⁵
Wild type Protein level	FYPO_EXT
Physical Interaction	BioGRID ⁶
Protein Modification	PSI-MOD ⁷

712 ¹ The Entry Summary section includes information on which gene is being displayed in the gene-centric
713 results page. The UniProtKB accession number is used to automatically retrieve the name and function of
714 the protein, plus any cross-referenced identifiers from Ensembl Genomes and NCBI GenBank. The
715 section also displays the PHI-base 5 gene identifier (PHIG) and any of the high-level terms
716 (Supplementary file 3) annotated to the gene.

717 ² Data from UniProtKB, NCBI Taxonomy and PubMed are automatically retrieved, while all other data are
718 manually curated.

719 ³ PHIPO is the Pathogen-Host Interaction Phenotype Ontology.

720 ⁴ GO is the Gene Ontology.

721 ⁵ FYPO_EXT is the Fission Yeast Phenotype Ontology Extension.

722 ⁶ BioGRID is the Biological General Repository for Interaction Datasets.

723 ⁷ PSI-MOD is the Human Proteome Organization (HUPO) Proteomics Standards Initiative (PSI) Protein
724 Modifications Ontology.

725
726

727 Figure supplements

728 **Figure 3 – figure supplement 1.** Canto entity relationship diagram.

729 Simplified UML class diagram showing the relations between entities (things of interest) in a
730 Canto curation session. The numbers on the connecting lines represent the cardinality of the
731 relation, meaning how many of one entity can be related to another entity: 0..n means 'zero or
732 more'; 1..n means 'one or more'. Lines with a hollow arrowhead indicate that the target entity (at
733 the head of the arrow) is a generalization of the source entity (at the tail of the arrow). Boxes
734 outlined in bold indicate new entities which were added to support curation in PHI-Canto.

735

736 **Figure 4 – figure supplement 1.** Alternative curation step workflow.

737 The flow diagram represents the PHI-Canto curation process from beginning to end in 5 steps. It
738 is an alternative representation to the image depicted in Figure 4. During step 2 of the workflow,
739 the curator chooses either the gene annotation or genotype / metagenotype annotation process.
740 Multiple annotations can be made using both annotation processes which can then be
741 submitted for review.

742

743 **Figure 4 - figure supplement 2.** What you need to curate a publication into PHI-Canto.

744 **Figure 4 - figure supplement 3.** Instructions on how to look up a UniProtKB ID.

745 **Figure 5 – figure supplement 1.** Resources relied upon by PHI-Canto.

746 Source code

747 **Source code 1.** Main configuration file for PHI-Canto.

748 This is the main configuration file for PHI-Canto. Much of the configuration is inherited from
749 Canto, the original curation application from which PHI-Canto is derived. Lines containing
750 custom configuration for PHI-Canto have been indicated with comments.

751

Appendix 1. How to use Annotation Extensions.

This file provides information on Annotation Extensions (AE) and how to use them in PHI-Canto to curate a standard selection of experiments (Table 2). The first section provides four examples of using AEs for curating metagenotypes with pathogen-host interaction phenotypes. The second section provides examples of curating metagenotypes using the gene-for-gene phenotype workflow, including using the AEs for gene-for-gene interactions and inverse gene-for-gene interactions. The third section of this file illustrates three examples of using AEs for curating single species phenotypes.

Further information on how to use PHI-Canto to make annotations can be found in PHI-Canto's user documentation, available at <https://canto.phi-base.org/docs/index>.

SECTION 1: Annotation Extensions for curating pathogen-host interaction phenotypes on metagenotypes

When creating and annotating metagenotypes, it is advisable to also create and annotate a wild type control metagenotype where possible. This enables a better understanding of annotations made to altered metagenotypes.

(Note: it is also possible to use several of the AEs in the table documenting single species phenotype AEs, e.g. *penetrance* and *affected protein*)

If you have a metagenotype phenotype recording 'unaffected pathogenicity' (corresponds to footnote 1 in Table 2)

AE summary table

AE name	Cardinality	Available terms
compared to control genotype	0, 1	Metagenotype identifier
extent of infectivity	0, 1	'unaffected pathogenicity'
host tissue affected	0, <i>n</i>	BRENDA Tissue Ontology term
outcome of interaction	0, 1	'disease present', 'disease absent'

Example publication: The RhIR quorum-sensing receptor controls *Pseudomonas aeruginosa* pathogenesis and biofilm development independently of its canonical homoserine lactone autoinducer. ([PMID:28715477](https://pubmed.ncbi.nlm.nih.gov/28715477/))

Pathogen-host interaction phenotype

Control metagenotype

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
rhII+[WT level] <i>P. aeruginosa</i> (PA14)	wild type <i>C. elegans</i> (N2)	PHIPO:0001069	death of host organism with pathogen	Cell growth assay	agar plates	Fig 6a	infects_tissue whole body , has_penetrance 50% , interaction_outcome disease present

Altered metagenotype

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
rhlIΔ <i>P. aeruginosa</i> (PA14)	wild type <i>C. elegans</i> (N2)	PHIPO:0001069	death of host organism with pathogen	Cell growth assay	agar plates	Fig 6a	infects_tissue whole body , infective_ability unaffected pathogenicity , has_penetrance 50% , compared_to_control rhlI+[WT level] <i>Pseudomonas aeruginosa</i> (PA14) / wild type <i>Caenorhabditis elegans</i> (N2) , interaction_outcome disease present

If you have a metagenotype phenotype recording 'altered pathogenicity or virulence' (corresponds to footnote 2 in Table 2)

AE summary table

AE name	Cardinality	Available terms
compared to control genotype	0, 1	Metagenotype identifier
extent of infectivity	0, 1	'loss of pathogenicity', 'reduced virulence', 'increased virulence'
host tissue affected	0, <i>n</i>	BRENDA Tissue Ontology term
outcome of interaction	0, 1	'disease present', 'disease absent'

Example publication: A conserved fungal glycosyltransferase facilitates pathogenesis of plants by enabling hyphal growth on solid surfaces ([PMID:29020037](#))

A training video is available for the curation of this publication at <https://youtu.be/44XGoi6ljgk?t=1738>

Pathogen-host interaction phenotype

Control metagenotype

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
GT2+[WT level] <i>Z. tritici</i> (IPO323)	wild type <i>T. aestivum</i> (cv. Riband)	PHIPO:0000480	presence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	14 days post inoculation	Figure 2E	infects_tissue leaf , interaction_outcome disease present

Altered metagenotype

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
ΔGT2-19(deletion) <i>Z. tritici</i> (IPO323)	wild type <i>T. aestivum</i> (cv. Riband)	PHIPO:0000481	absence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	14 days post inoculation	Figure 2E	infects_tissue leaf , infective_ability loss of pathogenicity , compared_to_control GT2+[WT level] <i>Zymoseptoria tritici</i> (IPO323) / wild type <i>Triticum aestivum</i> (cv. Riband) , interaction_outcome disease absent

If you have a metagenotype phenotype recording 'mutualism' (corresponds to footnote 4 in Table 2)

AE summary table

AE name	Cardinality	Available terms
---------	-------------	-----------------

compared to control genotype	0, 1	Metagenotype identifier
extent of infectivity	0, 1	'mutualism present', 'mutualism absent', 'loss of mutualism'
host tissue affected	0, <i>n</i>	BRENDA Tissue Ontology term

Note: The 'Outcome of interaction' AE is not relevant in this mutualism interaction.

Example publication: Reactive oxygen species play a role in regulating a fungus-perennial ryegrass mutualistic interaction ([PMID:16517760](#))

Pathogen-host interaction phenotype: Example 1

Illustrating a phenotype associated with the pathogen component within the Pathogen-Host Interaction.

Control metagenotype

Pathogen genotype	Host genotype	Term ID ↕	Term name ↕	Evidence code ↕	Figure ↕	Annotation extension ↕
noxA+[WT level] <i>E. festucae</i> (F1) bkg: GFP	wild type <i>L. perenne</i> (Unknown strain)	PHIPO:0000954	presence of pathogen growth within host	Microscopy	Figure 1c	infects_tissue leaf , infective_ability mutualism present

Altered metagenotype

Pathogen genotype	Host genotype	Term ID ↕	Term name ↕	Evidence code ↕	Figure ↕	Annotation extension ↕
noxA::pAN7-1(disruption)[Not assayed] <i>E. festucae</i> (F1) bkg: GFP	wild type <i>L. perenne</i> (Unknown strain)	PHIPO:0000368	increased pathogen growth within host	Microscopy	Figure 1d	infects_tissue leaf , infective_ability loss of mutualism , compared_to_control noxA+[WT level] <i>Epichloe festucae</i> (F1) / wild type <i>Lolium perenne</i> (Unknown strain)

Pathogen-host interaction phenotype: Example 2

Illustrating a phenotype associated with the host component within the Pathogen-Host Interaction.

Control metagenotype

Pathogen genotype	Host genotype	Term ID ↕	Term name ↕	Evidence code ↕	Figure ↕	Annotation extension ↕
noxA+[WT level] <i>E. festucae</i> (F1)	wild type <i>L. perenne</i> (Unknown strain)	PHIPO:0001005	normal host morphology during pathogen invasion	Macroscopic observation (qualitative observation)	Figure 1a, 5c	infects_tissue whole plant , infective_ability mutualism present

Altered metagenotype

Note: in this case, two separate annotations were made to the same metagenotype.

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Figure	Annotation extension
noxA::pAN7-1(disruption)[Not assayed] <i>E. festucae</i> (F11)	wild type <i>L. perenne</i> (Unknown strain)	PHIPO:0001130	stunted host growth during pathogen colonization	Macroscopic observation (qualitative observation)	Figure 1a	infects_tissue whole plant , infective_ability loss of mutualism , compared_to_control noxA+[WT level] <i>Epichloe festucae</i> (F11) / wild type <i>Lolium perenne</i> (Unknown strain)
noxA::pAN7-1(disruption)[Not assayed] <i>E. festucae</i> (F11)	wild type <i>L. perenne</i> (Unknown strain)	PHIPO:0001131	increased number of host side shoots during pathogen colonization	Macroscopic observation (qualitative observation)	Figure 1a	infects_tissue whole plant , infective_ability loss of mutualism , compared_to_control noxA+[WT level] <i>Epichloe festucae</i> (F11) / wild type <i>Lolium perenne</i> (Unknown strain)

If you have a metagenotype phenotype recording 'a pathogen effector' (corresponds to footnote 5 in Table 2)

If you have a biotrophic or necrotrophic plant pathogen effector which is involved in a gene-for-gene interaction, please see the AEs for the 'gene-for-gene interaction' or 'inverse gene-for-gene interaction' workflow (Section 2).

Annotate the pathogen effector with the GO Biological Process term 'effector-mediated modulation of host process by symbiont' ([GO:0140418](#)) or a descendant. If the GO Molecular Function term is known, then this can also be annotated and linked to the relevant GO effector term via an annotation extension.

AE summary table

AE name	Cardinality	Available terms
compared to control genotype	0, 1	Metagenotype identifier
extent of infectivity	0, 1	'unaffected pathogenicity', 'loss of pathogenicity', 'reduced virulence', 'increased virulence'
host tissue affected	0, <i>n</i>	BRENDA Tissue Ontology term
outcome of interaction	0, 1	'disease present', 'disease absent'

Example publication: An effector protein of the wheat stripe rust fungus targets chloroplasts and suppresses chloroplast function. ([PMID:31804478](#))

GO biological process

Species	Gene	Term ID	Term name	Evidence code	Figure
<i>P. striiformis</i>	PSTG_12806	GO:0052034	effector-mediated suppression of host pattern-triggered immunity	IDA	Figure 3a

Note: 'effector-mediated suppression of host pattern-triggered immunity' ([GO:0052034](#)) is a descendant term of 'effector-mediated modulation of host process by symbiont' ([GO:0140418](#)).

GO molecular function

Species	Gene	Term ID	Term name	Evidence code	With Figure	Annotation extension
<i>P. striiformis</i>	PSTG_12806	GO:0005515	protein binding	IPI	petC Figure 5	part_of effector-mediated suppression of host defenses
<i>P. striiformis</i>	PSTG_12806	GO:0004857	enzyme inhibitor activity	IPI	petC Figure 5	has_regulation_target petC , occurs_at host cell chloroplast , part_of effector-mediated suppression of host defenses

Please note that in the case of a physical interaction (protein–protein interaction) between the pathogen and host gene products (PSTG_12806 and PetC in the example above, respectively) this information can be curated using the Physical Interaction curation workflow, documented in https://canto.phi-base.org/docs/physical_interaction_annotation.

Pathogen-host interaction phenotype

Control metagenotype

In this case, there are no metagenotype control annotations. This is because it is not possible to create and annotate a metagenotype comprising of an empty vector control within the pathogen component of the metagenotype.

Altered metagenotype

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
Pst_12806ΔSP(1-23)[Not assayed] <i>P. striiformis</i> (f. sp. tritici strain CYR32)	wild type <i>N. benthamiana</i> (Unknown strain)	PHIPO:0001015	decreased level of host defense-induced callose deposition	Microscopy	delivery mechanism: agrobacterium, + PTI inducer flg22	Figure 3a, b	infects_tissue leaf , infective_ability increased virulence
Pst_12806ΔSP(1-23)[Not assayed] <i>P. striiformis</i> (f. sp. tritici strain CYR32)	wild type <i>N. benthamiana</i> (Unknown strain)	PHIPO:0001128	effector-mediated suppression of host PAMP-triggered immunity present	Microscopy	delivery mechanism: agrobacterium, + PTI inducer flg22	Figure 3a, b, c	infective_ability increased virulence

SECTION 2: Annotation Extensions for curating gene-for-gene phenotypes on metagenotypes

If you have a metagenotype phenotype recording 'a gene-for-gene interaction' (corresponds to footnote 6 in Table 2)

Annotate the pathogen effector with the GO Biological process term 'effector-mediated modulation of host process by symbiont' ([GO:0140418](#)) or a descendant. If the GO Molecular Function term is known, then this can also be annotated and linked to the relevant GO effector term via an annotation extension.

AE summary table

AE name	Cardinality	Available terms
compared to control genotype	0, 1	Metagenotype identifier
gene-for-gene phenotype	0, 1	'incompatible interaction, recognizable pathogen effector present, functional host resistance gene present'

		<p>‘incompatible interaction, recognizable pathogen effector present, gain of functional host resistance gene’</p> <p>‘incompatible interaction, gain of recognizable pathogen effector, gain of functional host resistance gene’</p> <p>‘incompatible interaction, gain of recognizable pathogen effector, functional host resistance gene present’</p> <p>‘compatible interaction, recognizable pathogen effector present, functional host resistance gene absent’</p> <p>‘compatible interaction, recognizable pathogen effector absent, functional host resistance gene present’</p> <p>‘compatible interaction, recognizable pathogen effector present, compromised host resistance gene’</p> <p>‘compatible interaction, recognizable pathogen effector absent, functional host resistance gene absent’</p> <p>‘compatible interaction, recognizable pathogen effector absent, compromised functional host resistance gene’</p>
--	--	---

		<p>'compatible interaction, compromised recognizable pathogen effector, functional host resistance gene present'</p> <p>'metagenotype outcome overcome by external condition'</p>
host tissue affected	0, n	BRENDA Tissue Ontology term

Example publication: Activation of an Arabidopsis resistance protein is specified by the in planta association of its leucine-rich repeat domain with the cognate oomycete effector. ([PMID:20601497](#)).

GO biological process

Species	Gene	Term ID	Term name	Evidence code
<i>H. arabidopsidis</i>	ATR1	GO:0140418	effector-mediated modulation of host process by symbiont	IMP

GO molecular function

Species	Gene	Term ID	Term name	Evidence code	With	Figure	Annotation extension
<i>H. arabidopsidis</i>	ATR1	GO:0005515	protein binding	IPI	RPP1	Figure 4, 5	part_of effector-mediated modulation of host process by symbiont

Gene-for-gene phenotype

Control metagenotypes

Incompatible control

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
ATR1-Δ51(1-51)[Not assayed] <i>H. arabidopsidis</i> (Maks9) bkg: Citrine tag	RPP1+[Not assayed] <i>A. thaliana</i> (ecotype Ws-0) bkg: HA tag	PHIPO:0000192	presence of host-defense induced lesion by host hypersensitive response	Macroscopic observation (qualitative observation)	delivery mechanism: agrobacterium, heterologous species tobacco	Figure 3a	infects_tissue leaf , gene_for_gene_interaction incompatible interaction, recognizable pathogen effector present, functional host resistance gene present

Compatible control

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
ATR1-Δ51(1-51)[Not assayed] <i>H. arabidopsidis</i> (Maks9) bkg: Citrine tag	RPP1+[Not assayed] <i>A. thaliana</i> (ecotype Nd-0) bkg: HA tag	PHIPO:0000182	absence of host-defense induced lesion by host hypersensitive response	Macroscopic observation (qualitative observation)	delivery mechanism: agrobacterium, heterologous species tobacco	Figure 3b, 8b	infects_tissue leaf , gene_for_gene_interaction compatible interaction, recognizable pathogen effector absent, functional host resistance gene present

Altered metagenotype (shift from compatible to incompatible interaction)

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
ATR1-Δ51-D191G(1-51, D191G)[Not assayed] <i>H. arabidopsidis</i> (Maks9) bkg: Citrine tag	RPP1+[Not assayed] <i>A. thaliana</i> (ecotype Nd-0) bkg: HA tag	PHIPO:0000192	presence of host-defense induced lesion by host hypersensitive response	Macroscopic observation (qualitative observation)	delivery mechanism: agrobacterium, heterologous species tobacco	Figure 8b	infects_tissue leaf , gene_for_gene_interaction incompatible interaction, gain of recognizable pathogen effector, functional host resistance gene present , compared_to_control ATR1-delta51(1-51)[Not assayed] Hyaloperonospora arabidopsidis (Maks9) / RPP1+[Not assayed] Arabidopsis thaliana (ecotype Nd-0)

If you have a metagenotype phenotype recording an inverse gene-for-gene interaction (corresponds to footnote 8 in Table 2)

Annotate the pathogen effector with the GO Biological process term 'effector-mediated modulation of host process by symbiont' ([GO:0140418](#)) or a descendant. If the GO Molecular Function term is known, then this can also be annotated and linked to the relevant GO effector term via an annotation extension.

AE summary table

AE name	Cardinality	Available terms
compared to control genotype	0, 1	Metagenotype identifier
inverse gene-for-gene phenotype	0, 1	<p>'compatible interaction, functional pathogen necrotrophic effector present, functional host susceptibility locus present'</p> <p>'compatible interaction, functional pathogen necrotrophic effector present, gain of functional host susceptibility locus'</p> <p>'compatible interaction, gain of functional pathogen necrotrophic effector, functional host susceptibility locus present'</p> <p>'incompatible interaction, functional pathogen necrotrophic effector present, functional host susceptibility locus absent'</p> <p>'incompatible</p>

		<p>interaction, functional pathogen necrotrophic effector absent, functional host susceptibility locus present'</p> <p>'incompatible interaction, functional pathogen necrotrophic effector present, functional host susceptibility locus compromised'</p> <p>'incompatible interaction, compromised functional pathogen necrotrophic effector, functional host susceptibility locus present'</p> <p>'incompatible interaction, gain of functional pathogen necrotrophic effector, functional host susceptibility locus compromised'</p> <p>'metagenotype outcome overcome by external condition'</p>
host tissue affected	0, <i>n</i>	BRENDA Tissue Ontology term

Example publication: The cysteine rich necrotrophic effector SnTox1 produced by *Stagonospora nodorum* triggers susceptibility of wheat lines harboring Snn1. ([PMID:22241993](#)).

GO biological process

Species ↕	Gene ↕	Term ID ↕	Term name ↕	Evidence code ↕
<i>P. nodorum</i>	Tox1	GO:0080185	effector-mediated induction of plant hypersensitive response by symbiont	EXP

GO molecular function

Species ↕	Gene ↕	Term ID ↕	Term name ↕	Evidence code ↕	Annotation extension ↕
<i>P. nodorum</i>	Tox1	GO:0140295	pathogen-derived receptor ligand activity	EXP	has_input Snn1 , part_of effector-mediated induction of plant hypersensitive response by symbiont

Gene-for-gene phenotype

Control metagenotypes

Compatible control

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
Tox1+[WT level] <i>P. nodorum</i> (SN15)	Snn1+[WT level] <i>T. aestivum</i> (cv. Chinese Spring)	PHIPO:0000480	presence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	delivery mechanism: culture infiltration	Figure 1	infests_tissue leaf , inverse_gene_for_gene compatible interaction, functional pathogen necrotrophic effector present, functional host susceptibility locus present

Incompatible control

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
Tox1-(no endogenous copy)[Not assayed] <i>P. nodorum</i> (Sn79-1087)	Snn1+[WT level] <i>T. aestivum</i> (cv. Chinese Spring)	PHIPO:0000481	absence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	delivery mechanism: pathogen spore inoculation	Figure 5b	infests_tissue leaf , inverse_gene_for_gene incompatible interaction, functional pathogen necrotrophic effector absent, functional host susceptibility locus present

Altered metagenotypes

Shift from compatible to incompatible interaction

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
Tox1+[WT level] <i>P. nodorum</i> (SN15)	Snn1-ems237(unknown)[Not assayed] <i>T. aestivum</i> (cv. Chinese Spring)	PHIPO:0000481	absence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	delivery mechanism: culture infiltration	Figure 1	infests_tissue leaf , compared_to_control Tox1+[WT level] <i>Parastagonospora nodorum</i> (SN15) / <i>Snn1+[WT level]</i> <i>Triticum aestivum</i> (Chinese Spring) , inverse_gene_for_gene incompatible interaction, functional pathogen necrotrophic effector present, functional host susceptibility locus compromised

Shift from incompatible to compatible interaction

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
+Sn15Tox1A1(transformant, no endogenous copy)[Not assayed] <i>P. nodorum</i> (Sn79-1087)	Snn1+[WT level] <i>T. aestivum</i> (cv. Chinese Spring)	PHIPO:0000480	presence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	delivery mechanism: pathogen spore inoculation	Figure 5b	infests_tissue leaf , compared_to_control Tox1-(no endogenous copy)[Not assayed] <i>Parastagonospora nodorum</i> (Sn79-1087) / <i>Snn1+[WT level]</i> <i>Triticum aestivum</i> (Chinese Spring) , inverse_gene_for_gene compatible interaction, gain of functional pathogen necrotrophic effector, functional host susceptibility locus present

No shift compared to control, still an incompatible interaction, despite alteration to both pathogen and host genotypes

Note: the AEs capture the detail of what has occurred within this pathogen-host interaction.

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
+Sn15Tox1A1(transformant, no endogenous copy)[Not assayed] <i>P. nodorum</i> (Sn79-1087)	Snn1-ems237(unknown)[Not assayed] <i>T. aestivum</i> (cv. Chinese Spring)	PHIPO:0000481	absence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	delivery mechanism: pathogen spore inoculation	Figure 5b	infects_tissue leaf , compared_to_control Tox1-(no endogenous copy)[Not assayed] Parastagonospora nodorum (Sn79-1087) / Snn1+[WT level] Triticum aestivum (Chinese Spring) , inverse_gene_for_gene incompatible interaction, gain of functional pathogen necrotrophic effector, functional host susceptibility locus compromised

SECTION 3: Annotation Extensions for curating single species phenotypes (pathogen phenotypes or host phenotypes)

AE summary table

AE name	Cardinality	Available terms
affected proteins	2	UniProtKB accession number
assayed RNA	0, 1	UniProtKB accession number
assayed protein	0, 1	UniProtKB accession number
penetrance	0, 1	qualitative terms ('high', 'medium', 'low', or 'complete') or a quantitative value (a percentage)
severity	0, 1	'high', 'medium', 'low', 'variable severity'
observed in organ	0, 1	BRENDA Tissue Ontology term

Example of an in vitro pathogen phenotype (corresponds to footnote 3 in Table 2)

Example publication: A conserved fungal glycosyltransferase facilitates pathogenesis of plants by enabling hyphal growth on solid surfaces. ([PMID:29020037](https://pubmed.ncbi.nlm.nih.gov/29020037/))

A training video is available for the curation of this publication at <https://youtu.be/44XGoi6ljgk?t=1738>

Pathogen phenotype

Species (strain)	Genes	Genotype (allele and expression)	Term ID	Term name	Evidence code	Conditions	Figure
<i>Z. tritici</i> (IPO323)	GT2	ΔGT2-19(deletion)	PHIPO:0001212	decreased hyphal growth	Cell growth assay	water medium, agar plates	Figure 2E

Please note that in this curation example, no AEs were required.

Example of an in vitro pathogen chemistry phenotype (corresponds to footnote 9 in Table 2)

Example publication: The T788G mutation in the *cyp51C* gene confers voriconazole resistance in *Aspergillus flavus* causing aspergillosis. ([PMID:22314539](#))

Pathogen phenotype

Species (strain)	Genes	Genotype (allele and expression)	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
<i>A. flavus</i> (NRRL 3357)	<i>cyp51c</i>	<i>cyp51C</i> -T788G(aaS240A)[WT level]	PHIPO:0000590	resistance to voriconazole	Cell growth assay	liquid culture, minimal medium, + voriconazole	Table 3 (footnote d), text page 2602	has_severity high
<i>A. flavus</i> (NRRL 3357)	<i>cyp51c</i>	<i>cyp51C</i> -T161C(aaM54T)[WT level]	PHIPO:0001219	normal growth on voriconazole	Cell growth assay	liquid culture, minimal medium, + voriconazole	text on page 2602	

Example of an in vitro host phenotype (corresponds to footnote 7 in Table 2)

Example publication: Activation of an Arabidopsis resistance protein is specified by the in planta association of its leucine-rich repeat domain with the cognate oomycete effector. ([PMID:20601497](#))

Host phenotype

Species (strain)	Genes	Background	Genotype (allele and expression)	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
<i>A. thaliana</i> (ecotype Ws-0)	RPP1	GFP	RPP1-TIR(266-1221)[Not assayed]	PHIPO:0000467	presence of effector-independent host hypersensitive response	Macroscopic observation (qualitative observation)	delivery mechanism: agrobacterium, heterologous species tobacco	Figure 7a, c	observed_organ leaf
<i>A. thaliana</i> (ecotype Ws-0)	RPP1	GFP	RPP1-TIRNBS(590-1221)[Not assayed]	PHIPO:0001180	absence of effector-independent host hypersensitive response	Macroscopic observation (qualitative observation)	delivery mechanism: agrobacterium, heterologous species tobacco	Figure 7a	observed_organ leaf
<i>A. thaliana</i> (ecotype Ws-0)	RPP1	GFP	RPP1-TIR E158A(266-1221, E158A)[Not assayed]	PHIPO:0001180	absence of effector-independent host hypersensitive response	Macroscopic observation (qualitative observation)	delivery mechanism: agrobacterium, heterologous species tobacco	Figure 7c	observed_organ leaf

Appendix 2. Worked example of a curation session.

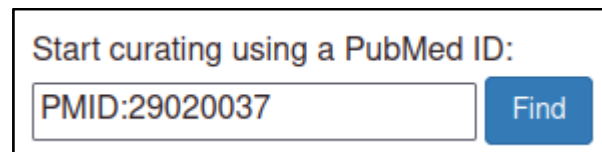
This document provides a worked example of the curation process in PHI-Canto for the publication by King et al. (2017), *A conserved fungal glycosyltransferase facilitates pathogenesis of plants by enabling hyphal growth on solid surfaces* ([PMID:29020037](#)).

The research study confirms the hypothesis that the GT2 gene is required for the fungal pathogens *Zymoseptoria tritici* and *Fusarium graminearum* to cause disease on wheat (*Triticum aestivum*). The curation session in PHI-Canto captures this conclusion by annotating a pathogen–host interaction between *Z. tritici* and *T. aestivum* to show that deletion of the GT2 gene causes loss of pathogenicity in the pathogen, and an absence of pathogen-associated lesions in the host. The wild type interaction between *Z. tritici* and *T. aestivum* is annotated to indicate the presence of disease (and lesions), and a corresponding pathogen–host interaction between *F. graminearum* and *T. aestivum* is annotated to show that deleting GT2 again causes a loss of pathogenicity and the absence of pathogen-associated lesions in the host.

The example starts with the entry of the publication into PHI-Canto (<https://canto.phi-base.org/>) and ends with the submission of the curation session for review by curators at PHI-base. The information curated from this publication is available on the new gene centric PHI-base 5 website (<http://phi5.phi-base.org>, search for PHIG:308 and PHIG:307).

Entering the publication

The PHI-Canto homepage provides a text field where publications can be entered by providing their PubMed ID (PMID). The PMID in this case is 29020037.



The image shows a web form for searching publications. At the top, it says "Start curating using a PubMed ID:". Below this is a text input field containing "PMID:29020037" and a blue button labeled "Find".

PHI-Canto will automatically retrieve details of the publication from PubMed so that the curator can confirm that they have entered the correct PMID.

Publication details

ID	PMID:29020037
Title	A conserved fungal glycosyltransferase facilitates pathogenesis of plants by enabling hyphal growth on solid surfaces.
Authors	King R, Urban M, Lauder RP, Hawkins N, Evans M, Plummer A, Halsey K, Lovegrove A, Hammond-Kosack K, Rudd JJ
Abstract	Pathogenic fungi must extend filamentous hyphae across solid surfaces to cause diseases of plants. However, the full inventory of genes which support this is incomplete and many may be currently concealed due to their essentiality for the hyphal growth form. During a random T-DNA mutagenesis screen performed on the pleomorphic wheat (<i>Triticum aestivum</i>) pathogen <i>Zymoseptoria tritici</i> , we acquired a mutant unable to extend hyphae specifically when on solid surfaces. In contrast "yeast-like"

After accepting the publication, the curator is prompted for their name, email address, and (optionally) an ORCID ID, which are used to attribute the curation to the curator, and to contact the curator in case of problems with the curation session.

Curator details

Before you start curating, please confirm your name and email address:

Name

Email

Your [ORCID](#) (optional but recommended):

[Why we collect ORCIDs](#)

Specifying genes and species

The gene is the most basic unit of annotation in PHI-Canto: every other biological feature that can be annotated involves a gene, so genes are entered first. PHI-Canto uses accession numbers from the UniProt Knowledgebase (UniProtKB) to uniquely identify proteins for the genes of interest in the curated publication.

The UniProtKB accession numbers for the publication are shown below.

Create gene list for PMID:29020037

Please list the genes studied in this paper using the UniProt identifier (eg. Q00909) separated by commas, spaces, tabs or one per line.

If you have large datasets please consider our [bulk annotation formats](#).

Note: Only supply high confidence interactions for large datasets.

You can edit this list later if you need to add more genes or remove "unused" genes.

F9WWD1 I1RB03

Since this publication describes a wild type host species (*T. aestivum*) with no specified genes of interest, the curator must add the host to the session by entering its NCBI Taxonomy ID in a separate field.

Add host organisms (where the paper has a host with no specified genes):

NCBI Taxon Id	Species	Common name (where available)	
4565	Triticum aestivum	bread wheat	X

PHI-Canto automatically retrieves details of the proteins from UniProtKB, including the gene name, gene product, and taxonomy (e.g., the species name).

Pathogen genes

Organism	Gene			
	ID	Name	Product	
<i>Fusarium graminearum</i>	I1RB03	GT2	Type 2 glycosyltransferase	X
<input type="text" value="Type a strain name"/>	<input type="button" value="Add strain"/>	<input type="button" value="Unknown strain"/>		
<i>Zymoseptoria tritici</i>	F9WWD1	GT2	Type 2 glycosyltransferase	X
<input type="text" value="Type a strain name"/>	<input type="button" value="Add strain"/>	<input type="button" value="Unknown strain"/>		

Host genes

Organism	Gene			
	ID	Name	Product	
<i>Triticum aestivum</i>	(No genes for this organism)			X
<input type="text" value="Type a strain name"/>	<input type="button" value="Add strain"/>	<input type="button" value="Unknown strain"/>		

Specifying strains

The curator must enter the strains for each organism studied in the publication or must specify when the strain was not known (or not specified in the publication). PHI-Canto provides a pre-populated list of strains for many species that the curator can select from, though they also have the option to specify a strain not in the list as free text.

In this publication, the pathogen strains are PH-1 for *F. graminearum* and IPO323 for *Z. tritici*. Two cultivars of *T. aestivum* were used: cv. Bobwhite and cv. Riband.

Pathogen genes

Organism	Gene			
	ID	Name	Product	
<i>Fusarium graminearum</i>	I1RB03	GT2	Type 2 glycosyltransferase	X
PH-1 X				
<input type="text" value="Type a strain name"/> <input type="button" value="Add strain"/> <input type="button" value="Unknown strain"/>				
<i>Zymoseptoria tritici</i>	F9WWD1	GT2	Type 2 glycosyltransferase	X
IPO323 X				
<input type="text" value="Type a strain name"/> <input type="button" value="Add strain"/> <input type="button" value="Unknown strain"/>				

Host genes

Organism	Gene			
	ID	Name	Product	
<i>Triticum aestivum</i>	(No genes for this organism)			X
cv. Bobwhite X cv. Riband X				
<input type="text" value="Type a strain name"/> <input type="button" value="Add strain"/> <input type="button" value="Unknown strain"/>				

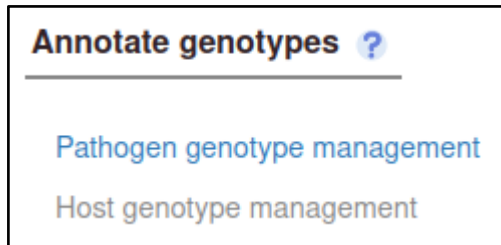
Creating alleles and genotypes

In order to show that deleting GT2 in the pathogen causes a loss of pathogenicity, the curator must annotate the interaction between the mutant pathogen and its host with a phenotype, meaning the interaction must be added to the curation session. In PHI-Canto, interactions are represented as *metagenotypes*, which are the combined genotypes of the pathogen and host species.

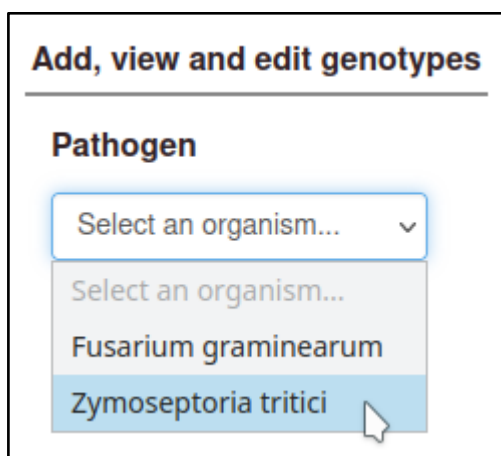
Before the curator can create a metagenotype, they must first create a genotype. Genotypes are composed from alleles (except in the case of wild type host genotypes with no specified genes, as described later), and metagenotypes are composed from genotypes. So, the

curator must first create an allele from a gene, then a genotype from an allele, then a metagenotype from two genotypes.

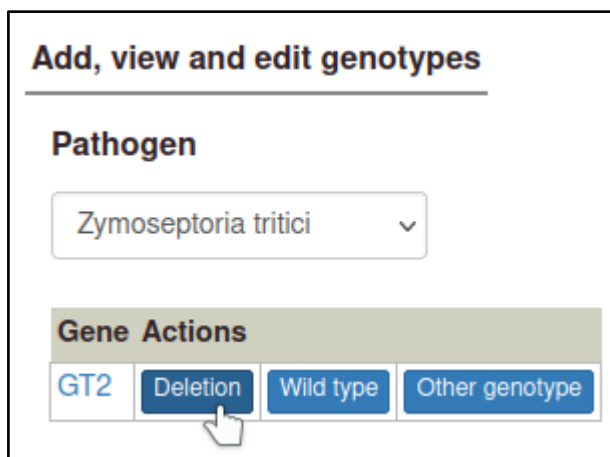
The curator starts from the Pathogen genotype management page, following a link from the Curation summary page.



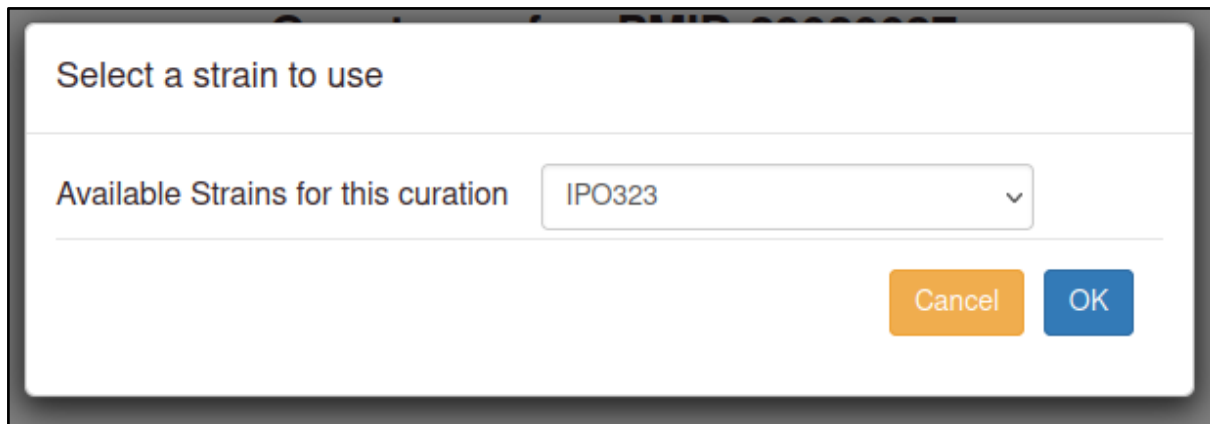
The curator then selects a pathogen species (*Z. tritici*) from a drop-down menu.



Selecting a pathogen species shows a list of genes for the species, with buttons to create types of alleles. Here, the curator selects 'Deletion' for a deletion allele.



The curator is prompted for the strain the deletion occurred in.

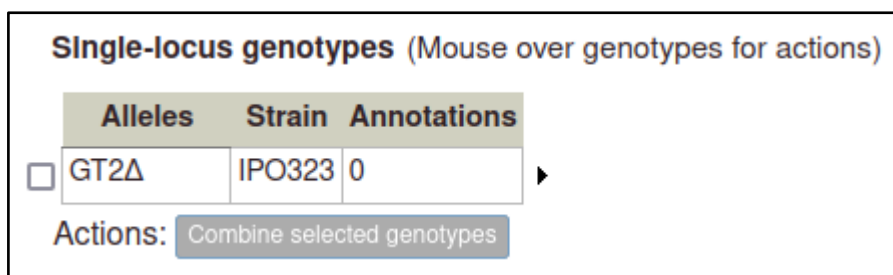


Select a strain to use

Available Strains for this curation IPO323

Cancel OK

After selecting this, PHI-Canto creates a genotype containing a single allele, with the allele name automatically generated from the gene name followed by a delta symbol.

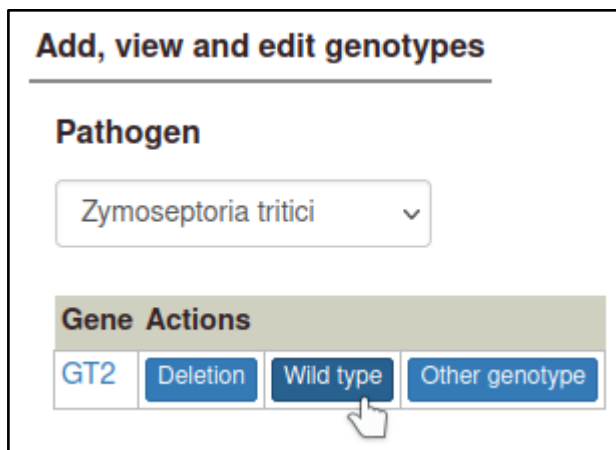


Single-locus genotypes (Mouse over genotypes for actions)

	Alleles	Strain	Annotations
<input type="checkbox"/>	GT2Δ	IPO323	0

Actions: Combine selected genotypes

The curator will also need to prepare a wild type genotype for the pathogen GT2 gene, which can be added to the control metagenotype so that any changes in the phenotype (between the wild type pathogen and the altered pathogen inoculated onto the host) can be properly annotated. This first requires making a wild type allele for GT2, using the 'Wild type' allele type.



Add, view and edit genotypes

Pathogen

Zymoseptoria tritici

Gene Actions

GT2 Deletion Wild type Other genotype

Wild type alleles require the gene expression level to be specified. In this case, there was no change in expression level, so the curator selects 'Wild type product level'. PHI-Canto automatically creates an allele name by appending a plus symbol to the gene name.

Adding allele for GT2

Allele name

Strain used

Expression ? Overexpression
 Wild type product level
 Knockdown
 Not assayed

As genotypes are created, they are added to a table of genotypes on their respective genotype management page (Pathogen genotype management for pathogens, Host genotype management for hosts).

Single-locus genotypes (Mouse over genotypes for actions)

	Alleles	Strain	Annotations	
<input type="checkbox"/>	GT2 Δ	IPO323	0	▶
<input type="checkbox"/>	GT2+[WT level]	IPO323	0	▶

Actions:

The curator can repeat the process above to create pathogen genotypes for *F. graminearum*.

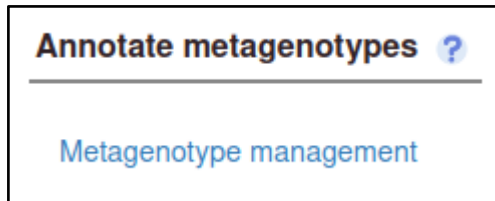
Single-locus genotypes (Mouse over genotypes for actions)

	Alleles	Strain	Annotations	
<input type="checkbox"/>	GT2 Δ	PH-1	0	▶
<input type="checkbox"/>	GT2+[WT level]	PH-1	0	▶

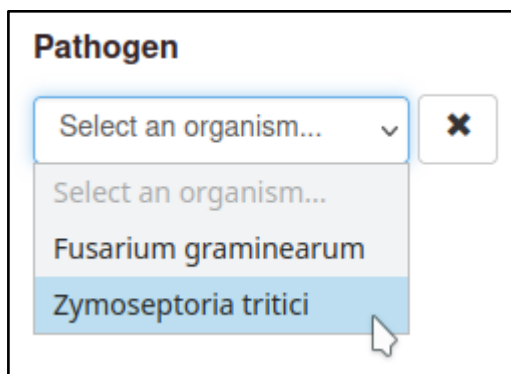
Actions:

Creating metagenotypes for pathogen–host interactions

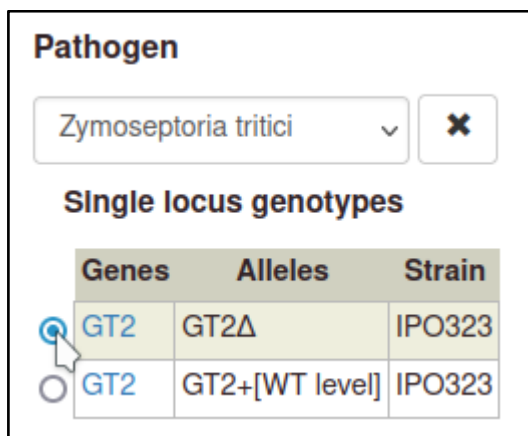
Metagenotypes are created using the Metagenotype management page, where genotypes previously added to the curation session can be combined into a metagenotype. The curator can reach this page from the Curation Summary page, or from either the pathogen or host genotype management page.



The curator starts by selecting a pathogen species from a drop-down menu.



Then the curator selects a genotype from the table of pathogen genotypes.



Then the curator selects a host genotype. For wild type hosts, PHI-Canto provides a shortcut where a strain can be selected without needing to create an allele as part of the genotype.

Host

Triticum aestivum

Wild type genotypes

cv. Bobwhite

cv. Riband

The curator selects 'Make metagenotype' to create the metagenotype for the interaction.

[← Go to Summary](#) [Make metagenotype](#)

The metagenotype is displayed in a table as a combination of pathogen and host genotype.

Metagenotypes

Pathogen		Host		Annotations
Species (strain)	Genotype	Species (strain)	Genotype	
Z. tritici (IPO323)	GT2Δ	T. aestivum (cv. Riband)	wild type	0

This process can be repeated to create the metagenotype for the wild type interaction between *Z. tritici* and *T. aestivum*. In this case, the pathogen genotype containing the wild type GT2 is selected instead of the deletion allele.

Pathogen

Zymoseptoria tritici

Single locus genotypes

	Genes	Alleles	Strain
<input type="radio"/>	GT2	GT2Δ	IPO323
<input checked="" type="radio"/>	GT2	GT2+[WT level]	IPO323

Metagenotypes

Pathogen		Host		
Species (strain)	Genotype	Species (strain)	Genotype	Annotations
<i>Z. tritici</i> (IPO323)	GT2Δ	<i>T. aestivum</i> (cv. Riband)	wild type	0
<i>Z. tritici</i> (IPO323)	GT2+[WT level]	<i>T. aestivum</i> (cv. Riband)	wild type	0

Creating the corresponding metagenotypes for *F. graminearum* and *T. aestivum* simply requires changing the pathogen species and selecting cv. Bobwhite for the host strain.

Metagenotypes

Pathogen		Host		
Species (strain)	Genotype	Species (strain)	Genotype	Annotations
<i>F. graminearum</i> (PH-1)	GT2Δ	<i>T. aestivum</i> (cv. Bobwhite)	wild type	0
<i>F. graminearum</i> (PH-1)	GT2+[WT level]	<i>T. aestivum</i> (cv. Bobwhite)	wild type	0

Annotating pathogen–host interactions with phenotypes

Metagenotypes can be annotated with phenotypes by selecting the ‘Annotate pathogen-host interaction phenotype’ action.

Metagenotypes

Pathogen		Host		Annotations
Species (strain)	Genotype	Species (strain)	Genotype	
Z. tritici (IPO323)	GT2Δ	T. aestivum (cv. Riband)	wild type	0

[Annotate pathogen-host interaction phenotype](#)
[Annotate gene-for-gene phenotype](#)
[Annotate disease name](#)
[View phenotype annotations](#)
[Delete](#)

Phenotype and evidence

The first step is to select a term from a controlled vocabulary that describes the phenotype of the interaction. PHI-Canto uses terms from the Pathogen–Host Interaction Phenotype Ontology (PHIPO) for this purpose. The primary observed phenotype in this case is the *absence of pathogen-associated host lesions* (PHIPO:0000481).

Annotating metagenotype – GT2delta *Zymoseptoria tritici* (IPO323) / wild type *Triticum aestivum* (cv. Riband)

Search for pathogen-host interaction phenotype term

Annotate normal or abnormal phenotypes of organisms within this pathogen-host interaction (metagenotype).

[more...](#)

Start typing a PHI phenotype in the search box (type at least 2 characters). If you do not find the term you are looking for with your initial search, begin with a broad term (pathogen colonization of host phenotype, binding, effector, host lesion) [more...](#)

absence lesions

- absence of pathogen-associated host lesions (PHIPO:0000481)
- presence of pathogen-associated host lesions (PHIPO:0000480)
- decreased extent of pathogen-associated host lesions (PHIPO:0000985)
- increased extent of pathogen-associated host lesions (PHIPO:0000986)
- presence of pathogen-associated host defense induced lesions (PHIPO:0000461)
- absence of pathogen growth within host (PHIPO:0000363)
- absence of pathogen growth on host surface (PHIPO:0000350)

Term name

absence of pathogen-associated host lesions

Definition

A phenotype where the process of host tissue cell death causing a host lesion is absent.

Upon selecting the term, the curator is shown a description of the term and its synonyms to help confirm that their chosen term is appropriate.

Annotating metagenotype – GT2delta *Zymoseptoria tritici* (IPO323) / wild type *Triticum aestivum* (cv. Riband)

Please read the term definition to ensure that it accurately describes your metagenotype

ID	PHIPO:0000481
Ontology	pathogen_host_interaction_phenotype
Term name	absence of pathogen-associated host lesions
Definition	A phenotype where the process of host tissue cell death causing a host lesion is absent.
Comment	The lesion can be induced by either the pathogen directly killing host tissue (e.g. cell wall degradation), or the host activating its own cell death pathways in defense. Note that if you are curating a necrotroph you need to annotate to PHIPO:0000465.
Synonyms	

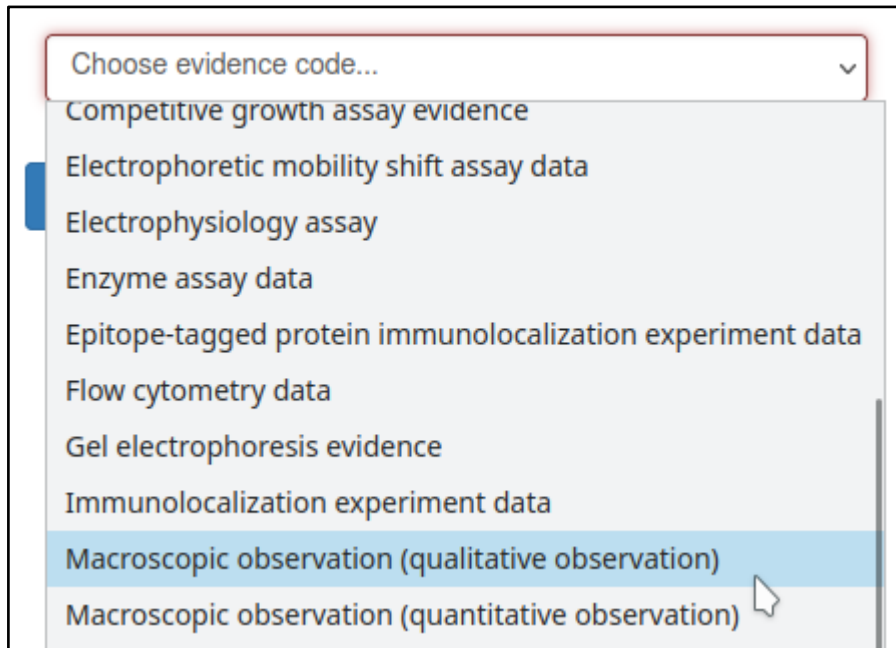
Can you use a more specific available term?

- [absence of host-defense induced lesion by host hypersensitive response](#) →
- [absence of pathogen necrotrophic effector-mediated host programmed cell death](#) →

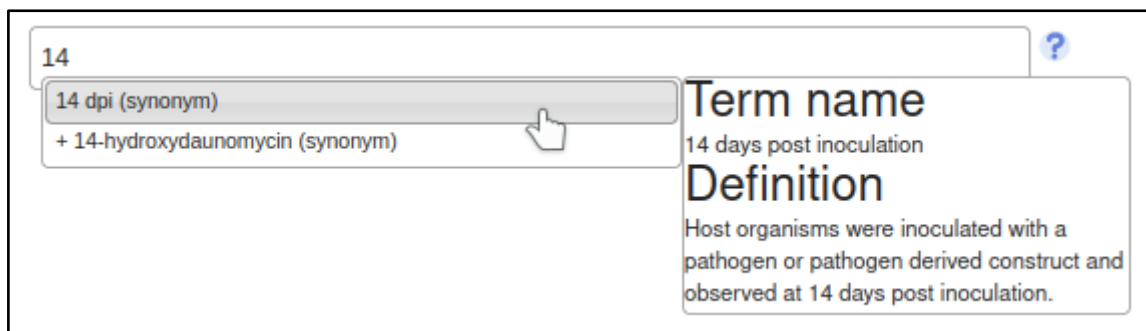
If you need a more specific term to describe the experiment you are annotating, and if none of terms above is appropriate, you can suggest a new term:

[Suggest a new child term for PHIPO:0000481](#)

The curator must select an evidence code for the observation of the phenotype. In this case, the phenotype was observed macroscopically, and measured qualitatively.



The curator may also specify experimental conditions for the experiment – such as the growth medium, or days elapsed after inoculation of the host. This annotation specifies that the assay was performed 14 days after inoculation with the *Z. tritici* GT2 deletion mutant.



Annotation extensions

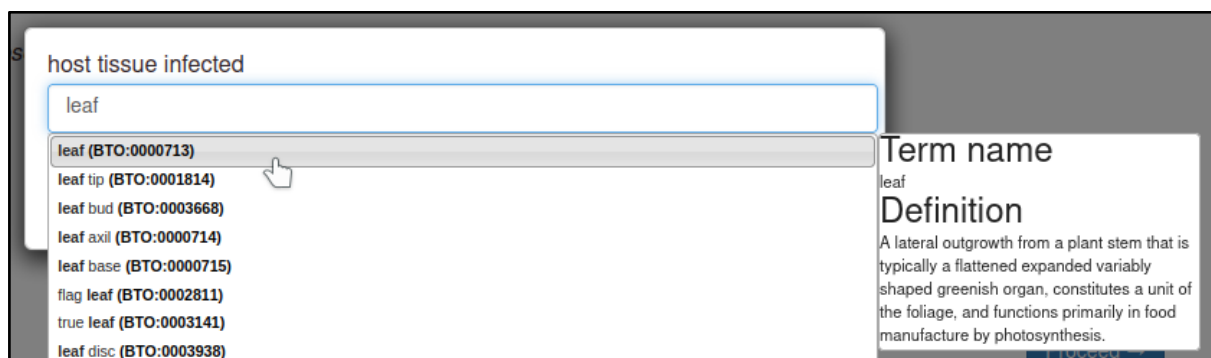
PHI-Canto uses annotation extensions to provide additional information about the conditions and outcome of the pathogen–host interaction. Of particular note are the host tissue infected, the changes to the infective ability of the pathogen, the presence (or absence) of disease, and the interaction used as a control for the interaction involving a mutant pathogen.

Annotation extensions

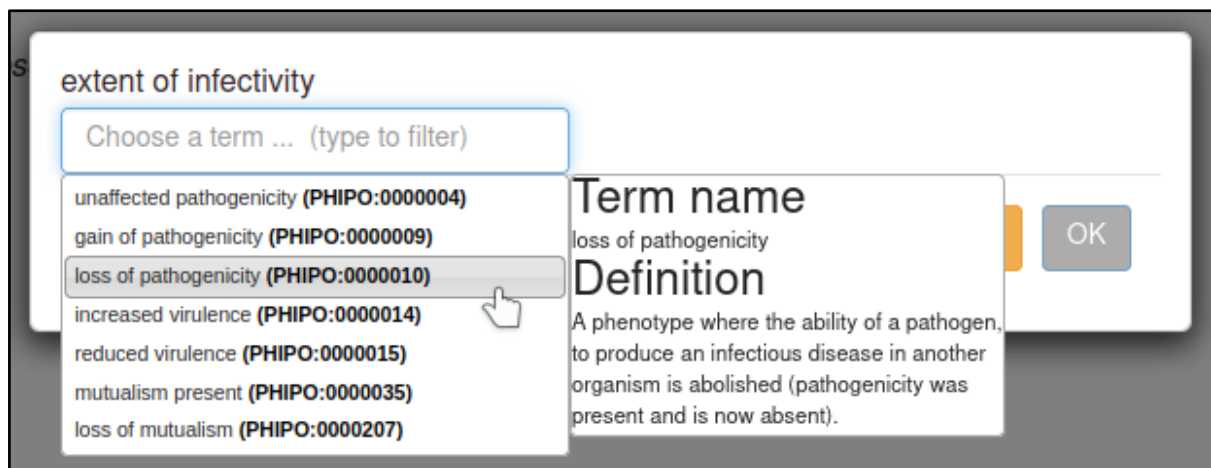
These extension types are available for *absence of pathogen-associated host lesions* (PHIPO:0000481):

- compared to control genotype
- penetrance
- severity
- extent of infectivity
- host tissue infected
- outcome of interaction

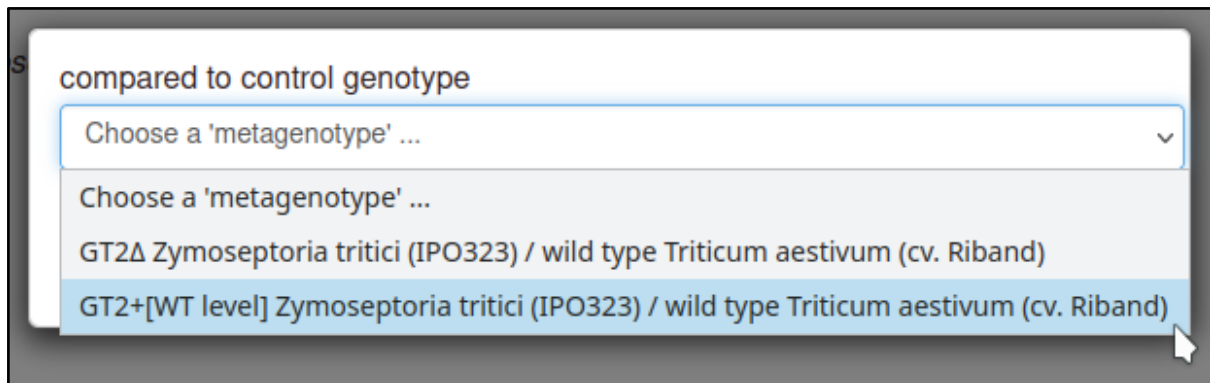
The host tissue that was infected during the interaction is annotated with the 'host tissue infected' annotation extension. This extension uses ontology terms from the BRENDA Tissue Ontology (BTO). In this case, the curator specifies that the *leaf* (BTO:0000713) of *T. aestivum* was infected.



Changes in the infective ability of the pathogen are annotated with the 'extent of infectivity' annotation extension. This extension uses a subset of ontology terms from PHIPO. In this case, the curator specifies that the interaction resulted in a *loss of pathogenicity* (PHIPO:0000010).

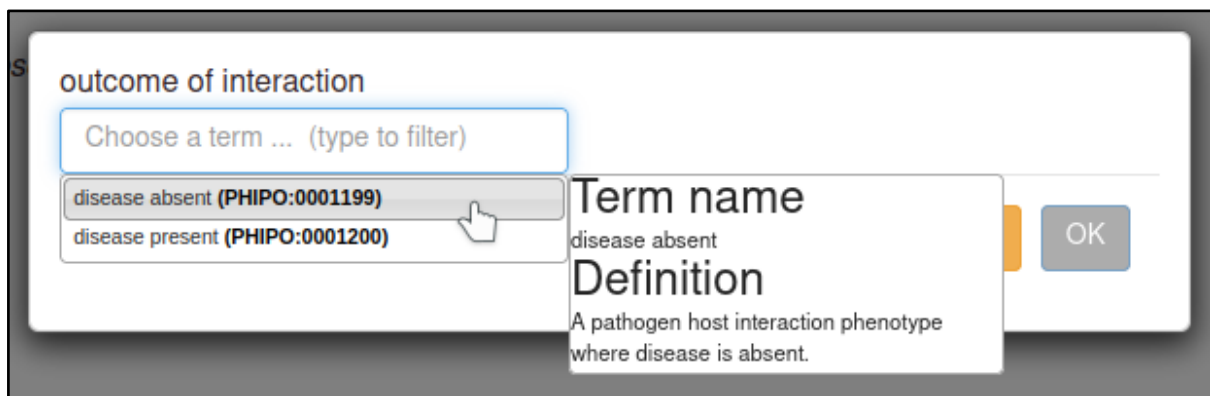


The control interaction (to which the interaction being annotated should be compared) can be annotated with the 'compared to control genotype' annotation extension. This annotation allows any metagenotype in the curation session to be designated as a control. In this case, the curator selects the wild type metagenotype that was created earlier.



The screenshot shows a dropdown menu titled "compared to control genotype". The menu is open, displaying a search bar with the text "Choose a 'metagenotype' ..." and a list of options. The selected option is "GT2+[WT level] Zymoseptoria tritici (IPO323) / wild type Triticum aestivum (cv. Riband)".

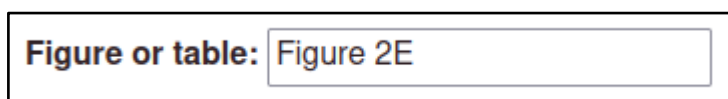
The presence or absence of disease resulting from the interaction can be annotated with the 'outcome of interaction' annotation extension. This extension uses a subset of ontology terms from PHIPO. In this case, the curator specifies that no disease was observed as a result of the interaction: *disease absent* (PHIPO:0001199).



The screenshot shows the "outcome of interaction" selection interface. A search bar contains the text "Choose a term ... (type to filter)". Below the search bar, two terms are listed: "disease absent (PHIPO:0001199)" and "disease present (PHIPO:0001200)". A mouse cursor is hovering over the "disease absent" term. To the right, a tooltip displays the term name "disease absent" and its definition: "A pathogen host interaction phenotype where disease is absent." An "OK" button is visible in the bottom right corner of the tooltip.

Figure numbers and comments

After adding annotation extensions, the curator has the option to provide the figure number from the publication (if any) that illustrates the phenotype. In this case, the figure was Figure 2E.



The screenshot shows a text input field labeled "Figure or table:" with the text "Figure 2E" entered.

The curator can also provide additional information in a comments field, in case of details that are not appropriate for any other field.

Once the above steps are completed, the phenotype annotation is created.

Pathogen-host Interaction phenotype							
Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
GT2Δ <i>Z. tritici</i> (IPO323)	wild type <i>T. aestivum</i> (cv. Riband)	PHIPO:0000481	absence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	14 days post inoculation	Figure 2E	infects_tissue leaf , infective_ability loss of pathogenicity , compared_to_control GT2+[WT level] Zymoseptoria tritici (IPO323) / wild type Triticum aestivum (cv. Riband) , interaction_outcome disease absent

Copying annotations

The above annotation can be used as a template for the interaction between the wild type pathogen and host, since many of the variables are the same. PHI-Canto provides a 'Copy and edit' feature that allows curators to use one annotation as a template for creating another.

Conditions	Figure	Annotation extension	
14 days post inoculation	Figure 2E	infects_tissue leaf , infective_ability loss of pathogenicity , compared_to_control GT2+[WT level]	View metagenotype Edit Copy and edit Delete

For the wild type interaction, the pathogen genotype is changed to wild type GT2, the phenotype term is changed to *presence of pathogen-associated host lesions* (PHIPO:0000480), the interaction outcome is changed to *disease present* (PHIPO:0001200), and the extensions for infective ability and control metagenotypes are removed, since they are not applicable.

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
GT2+[WT level] <i>Z. tritici</i> (IPO323)	wild type <i>T. aestivum</i> (cv. Riband)	PHIPO:0000480	presence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	14 days post inoculation	Figure 2E	infects_tissue leaf , interaction_outcome disease present

The interaction between *Z. tritici* and *T. aestivum* can also be used as a template for the interaction between *F. graminearum* and *T. aestivum*. Here, the pathogen genotype is changed to the GT2 deletion *F. graminearum*, the host strain is changed to cv. Bobwhite, the experimental condition is changed to '13 days post inoculation', the host tissue infected is changed to *inflorescence* (BTO:0000628), the control metagenotype is updated accordingly, and the figure number is changed to 4E.

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
GT2Δ <i>F. graminearum</i> (PH-1)	wild type <i>T. aestivum</i> (cv. Bobwhite)	PHIPO:0000481	absence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	13 days post inoculation	Figure 4E	infects_tissue inflorescence , infective_ability loss of pathogenicity , compared_to_control GT2+[WT level] Fusarium graminearum (PH-1) / wild type Triticum aestivum (cv. Bobwhite) , interaction_outcome disease absent

The changes required for the wild type interaction between *F. graminearum* and *T. aestivum* are the same as those required for *Z. tritici* and *T. aestivum*, since the interaction outcome is the same (presence of pathogen-associated host lesions, and presence of disease).

Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
GT2+[WT level] <i>F. graminearum</i> (PH-1)	wild type <i>T. aestivum</i> (cv. Bobwhite)	PHIPO:0000480	presence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	13 days post inoculation	Figure 4E	infects_tissue leaf , interaction_outcome disease present

Shown below is a table of all the pathogen–host interaction phenotypes from this curation example.

Pathogen-host Interaction phenotype							
Pathogen genotype	Host genotype	Term ID	Term name	Evidence code	Conditions	Figure	Annotation extension
GT2Δ <i>Z. tritici</i> (IPO323)	wild type <i>T. aestivum</i> (cv. Riband)	PHIPO:0000481	absence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	14 days post inoculation	Figure 2E	infects_tissue leaf , infective_ability loss of pathogenicity , compared_to_control GT2+[WT level] Zymoseptoria tritici (IPO323) / wild type Triticum aestivum (cv. Riband) , interaction_outcome disease absent
GT2+[WT level] <i>Z. tritici</i> (IPO323)	wild type <i>T. aestivum</i> (cv. Riband)	PHIPO:0000480	presence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	14 days post inoculation	Figure 2E	infects_tissue leaf , interaction_outcome disease present
GT2Δ <i>F. graminearum</i> (PH-1)	wild type <i>T. aestivum</i> (cv. Bobwhite)	PHIPO:0000481	absence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	13 days post inoculation	Figure 4E	infects_tissue inflorescence , infective_ability loss of pathogenicity , compared_to_control GT2+[WT level] Fusarium graminearum (PH-1) / wild type Triticum aestivum (cv. Bobwhite) , interaction_outcome disease absent
GT2+[WT level] <i>F. graminearum</i> (PH-1)	wild type <i>T. aestivum</i> (cv. Bobwhite)	PHIPO:0000480	presence of pathogen-associated host lesions	Macroscopic observation (qualitative observation)	13 days post inoculation	Figure 4E	infects_tissue leaf , interaction_outcome disease present

Disease annotation

PHI-Canto provides the 'Disease name' annotation type, which is used to annotate a disease to a pathogen–host interaction. These annotations highlight the fact that two different pathogens infecting different tissue types of the same host have been used in experiments within this publication.

Disease name annotations are made on the Metagenotype Management page, via the 'Annotate disease name' link.

Pathogen		Host		
Species (strain)	Genotype	Species (strain)	Genotype	Annotations
Z. tritici (IPO323)	GT2+[WT level]	T. aestivum (cv. Riband)	wild type	1

[Annotate pathogen-host interaction phenotype](#)
[Annotate gene-for-gene phenotype](#)
[Annotate disease name](#)
[View phenotype annotations](#)
[Delete](#)

The curator can select a disease from a list of disease names provided by the PHI-base Disease List (PHIDO). For *Z. tritici*, the disease is *septoria leaf blotch* (PHIDO:0000329).

septoria

- septoria leaf blotch (PHIDO:0000329)
- septoria nodorum blotch (PHIDO:0000330)
- septoria tritici blotch (PHIDO:0000331)

Term name
septoria leaf blotch

Definition
[no definition]

Disease name annotations also allow the host tissue infected to be specified. In this case, the tissue is the *leaf* (BTO:0000713).

host tissue infected

leaf

- leaf (BTO:0000713)
- leaf tip (BTO:0001814)
- leaf bud (BTO:0003668)
- leaf axil (BTO:0000714)
- leaf base (BTO:0000715)
- flag leaf (BTO:0002811)
- true leaf (BTO:0003141)
- leaf disc (BTO:0003938)

Term name
leaf

Definition
A lateral outgrowth from a plant stem that is typically a flattened expanded variably shaped greenish organ, constitutes a unit of the foliage, and functions primarily in food manufacture by photosynthesis.

The curator has the option to provide the figure number and additional comments. In this case, the figure numbers are 1 and 2.

Figure or table: Figure 1, 2

Once this step is completed, the disease name annotation is created.

Disease name					
Pathogen genotype	Host genotype	Term ID	Term name	Figure	Annotation extension
GT2+[WT level] <i>Z. tritici</i> (IPO323)	wild type <i>T. aestivum</i> (cv. Riband)	PHIDO:0000329	septoria leaf blotch	Figure 1, 2	infects_tissue leaf

The same process can be followed to create the Disease name annotation for *F. graminearum*: the genotype is the wild type GT2, the host cultivar is cv. *Bobwhite*, the disease is *fusarium ear blight* (PHIDO:0000162), the host tissue infected is the *inflorescence* (BTO:0000628), and the figure number is 4.

Disease name					
Pathogen genotype	Host genotype	Term ID	Term name	Figure	Annotation extension
GT2+[WT level] <i>F. graminearum</i> (PH-1)	wild type <i>T. aestivum</i> (cv. Bobwhite)	PHIDO:0000162	fusarium ear blight	Figure 4	infects_tissue inflorescence

Gene Ontology annotation

PHI-Canto also provides the ability to annotate biological processes, molecular functions, and cellular components associated with wild type versions of genes, using terms from the Gene Ontology (GO). In this publication, GT2 is described as having glycotransferase activity as its molecular function, so the curator can annotate this.

Gene Ontology annotations are made by selecting the gene from the Curation Summary page.

Annotate genes ?

Pathogens

- Fusarium graminearum
- GT2
- Zymoseptoria tritici
- GT2

The gene details page has a list of available annotation types.

Choose curation type for GT2: ?

- GO molecular function ?
- GO biological process ?
- GO cellular component ?
- Protein modification ?
- Physical interaction ?
- Wild-type RNA level ?
- Wild-type protein level ?
- Single allele phenotype ?

The curator selects the GO Molecular Function annotation type and is prompted for a term from the Gene Ontology. In this case, the correct term is *glycotransferase activity* (GO:0016757).

<input type="text" value="glycosyltransferase"/>	
glycosyltransferase activity (GO:0016757)	Term name
UDP-glycosyltransferase activity (GO:0008194)	glycosyltransferase activity
phenanthrol glycosyltransferase activity (GO:0019112)	Definition
peptidoglycan glycosyltransferase activity (GO:0008955)	Catalysis of the transfer of a glycosyl group from one compound (donor) to another (acceptor).
transfer ribonucleate glycosyltransferase activity (GO:0008479) (synonym)	
kinetin UDP glycosyltransferase activity (GO:0102694)	

The curator must provide an evidence code from a controlled list specified by the Gene Ontology. The appropriate evidence code in this case is a *Traceable Author Statement* in the publication.

Choose evidence for annotating GT2 with GO:0016757

Choose evidence code... ▼

- Choose evidence code...
- Inferred from Direct Assay (IDA)
- Inferred from Genetic Interaction (IGI)
- Inferred from Mutant Phenotype (IMP)
- Inferred from Physical Interaction (IPI)
- Inferred from Experiment (EXP)
- Traceable Author Statement (TAS)**

There are many annotation extensions available for GO annotations, but in this case, none of them are applicable (or required), so the curator skips this step.

Annotation extensions

These extension types are available for *glycosyltransferase activity* (GO:0016757):

with host species
has function during
physical location
involved in biological process
PR:nnn ID for gene product form
qualifier

Figure numbers can be specified for GO annotations: in this case, the relevant figure is Figure 3.

Figure or table:

Once this step is completed, the molecular function annotation is created.

GO molecular function

Species	Gene	Term ID	Term name	Evidence code	Figure
<i>Z. tritici</i>	GT2	GO:0016757	glycosyltransferase activity	TAS	Figure 3

Other annotation types

The publication contains other information which is not included in this worked example for the sake of brevity. In the real curation session, this other information is captured as the following annotations:

- **GO biological process** annotations indicate that GT2 is involved in the hyphal growth process.
- **GO cellular component** annotations indicate that GT2 is located in the hyphal cell wall.
- **Pathogen phenotype** annotations capture information about the pathogen *in vitro*, specifically normal and altered phenotypes for unicellular population growth, hyphal growth, cellular melanin accumulation, filament morphology, and so on.

All these annotation types use the same annotation process as the annotation types described above.

Submitting the curation session

Once the curator has made all their annotations, the curation session is submitted to the PHI-base team for review.



The curator can use a text box to provide any information that is outside the scope of the curation process before finishing the submission process. Once the submission process is finished, the curation session can no longer be edited except by members of the PHI-base team, who have the option to reactivate the session in case changes are required by the original curator.

Appendix 3. Author checklist prior to publication.

1. Use the most current gene name. Take care with synonyms. Prefix the gene name with the genus and species initials if the same gene name exists in multiple species.
2. If reporting on a new (gene) sequence, submit your sequence to NCBI GenBank or the European Nucleotide Archive (ENA), then obtain an accession number prior to publication. Record this accession number within the manuscript. If reporting on a gene with an existing accession number, make sure this is reported in the manuscript. Please record the UniProtKB accession number for the protein of the gene, where available. Provide or use any existing informative allele or line designations for mutations and transgenes.
3. Provide a binomial species name for pathogen and host organisms, not just a common name. If possible, please also include NCBI Taxonomy IDs for the pathogen and host organisms at the rank of species.
4. Describe the tissue or organ in which the experimental observations were made (controlled language can be found in the BRENDA Tissue Ontology, see <https://www.ebi.ac.uk/ols/ontologies/bto>).
5. Describe any experimental techniques used, and accurately record any chemicals or reagents used.
6. When writing an article, try to keep the use of descriptive language as accurate and controlled as possible. For example, do not use 'reduced pathogenicity' or 'loss of virulence', as these terms can be misleading: it would be more accurate to use 'reduced virulence' and 'loss of pathogenicity', respectively. Ideally, try to follow the terminology of an existing ontology: this will make the data easier to extract and reuse. Relevant ontologies include PHIPO and GO (<https://www.ebi.ac.uk/ols/ontologies/phipo>, <https://www.ebi.ac.uk/ols/ontologies/go>).
7. Document all the key information for the paper: do not rely on citing past papers for information on the pathogen used, or the strain used, and so on.
8. Think carefully when choosing keywords for your manuscript to ensure that the publication can be located by PHI-base's keyword searches. One example of an ideal keyword is 'pathogen-host interaction'.
9. Record the provenance of the pathogen strain: for example, whether it is a lab strain or a field isolate, or if the strain was obtained from a stock center or as a gift from another lab.

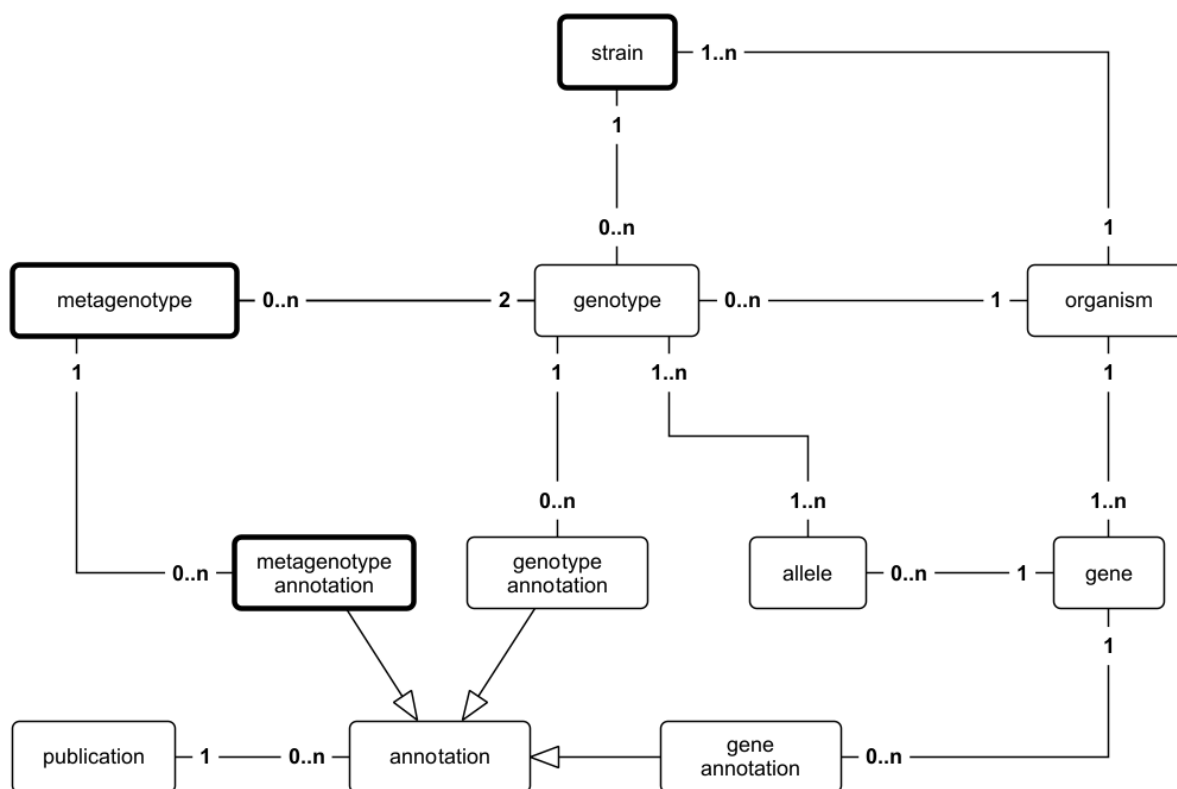


Figure 3 – figure supplement 1. Canto entity relationship diagram.

Simplified UML class diagram showing the relations between entities (things of interest) in a Canto curation session. The numbers on the connecting lines represent the cardinality of the relation, meaning how many of one entity can be related to another entity: 0..n means ‘zero or more’; 1..n means ‘one or more’. Lines with a hollow arrowhead indicate that the target entity (at the head of the arrow) is a generalization of the source entity (at the tail of the arrow). Boxes outlined in bold indicate new entities which were added to support curation in PHI-Canto.

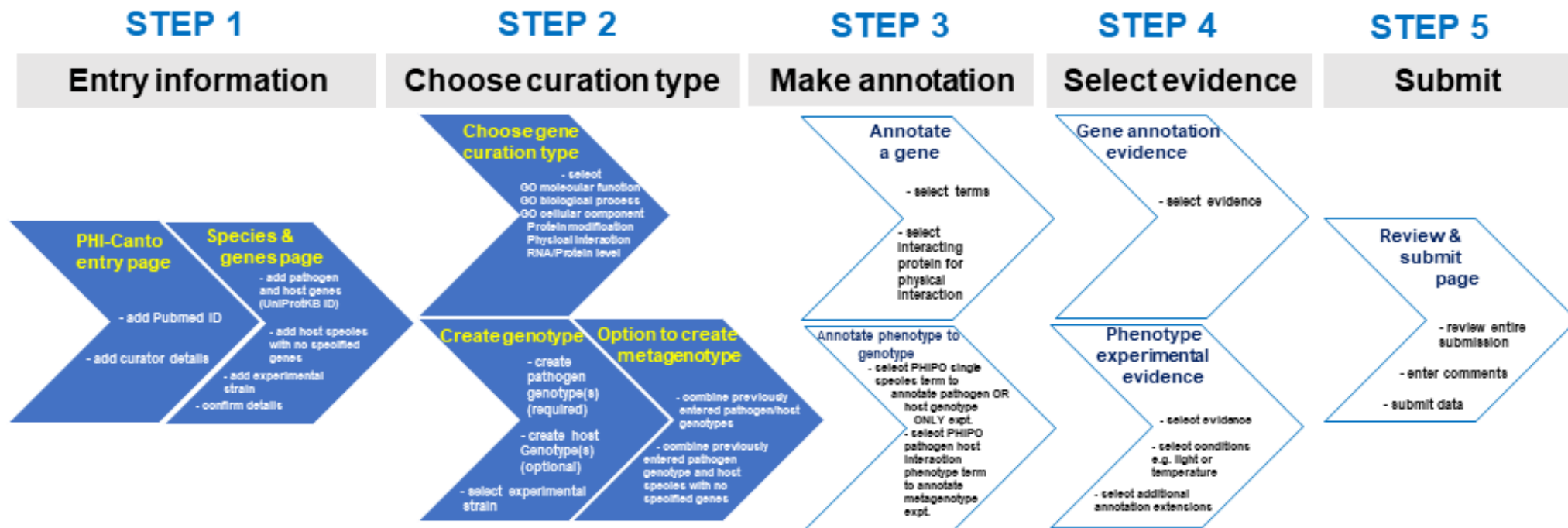


Figure 4 – figure supplement 1. Alternative curation step workflow.

The flow diagram represents the PHI-Canto curation process from beginning to end in 5 steps. It is an alternative representation to the image depicted in Figure 4. During step 2 of the workflow, the curator chooses either the gene annotation or genotype / metagenotype annotation process. Multiple annotations can be made using both annotation processes which can then be submitted for review.

Figure 4 - figure supplement 2. What you need to curate a publication into PHI-Canto.

1. The PubMed ID of the peer-reviewed publication.
2. Your email address, so we can contact you regarding your curation session.
3. UniProtKB accession numbers for the pathogen and host gene products studied within the publication.
4. The binomial names of the pathogen and host species studied within the publication.
5. Details of the experimental strains used within the publication.

Figure 4 - figure supplement 3. Instructions on how to look up a UniProtKB ID.

UniProtKB is divided into two sources: UniProtKB/Swiss-Prot, which contains manually annotated entries; and UniProtKB/TrEMBL, which contains unreviewed entries that are automatically annotated by prediction systems. PHI-Canto permits annotations on entries from either source, although UniProtKB/Swiss-Prot entries are preferred (owing to their higher quality).

Finding genes in UniProtKB

PHI-Canto uses [UniProt Knowledgebase](#) (UniProtKB) gene accession numbers to disambiguate genes/proteins. This is to ensure that we are talking about the correct gene product – especially as the same names are sometimes used for different proteins – and to standardize entries, because not all strains are in UniProt.

1. **Identify the reference proteome** (we use the designated reference proteome to integrate different strain information at the gene level in PHI-base). In PHI-Canto you will be able to specify the strain you used.

Look up the reference proteome for your organism using the species name (https://www.uniprot.org/help/reference_proteome).

If there is no reference proteome use the strain studied.

2. **Identify the gene of interest** in the reference proteome:

Start from the [UniProt homepage](#), then perform any of the following steps:

Search for the author assigned gene name/primary name (e.g., Tri5) or synonyms, plus species name (e.g., *Fusarium graminearum*).

OR

If the gene does not have a 'given name' but a locus ID is provided, search using the locus_id (e.g., FGRRES_03537) plus species name (e.g., *Fusarium graminearum*). If the entry identifier used is not the reference strain, copy the protein sequence and go to the BLAST step below.

OR

Search on a protein description (e.g., **Trichodiene synthase**)

OR

Obtain the protein sequence for your gene of interest and BLAST against UniProtKB (<https://www.uniprot.org/blast/>) with your protein sequence.

Note: If there are multiple entries for your gene product from the reference strain, please select the 'Reviewed entry'. Use the left-hand filter for 'Reviewed entries'.

OR

If the gene cannot be located in UniProt, contact the authors, UniProt, or PHI-base for help locating the canonical database entry.

- 3. Add the entry into PHI-Canto.** Once the entry of interest is located, select the entry accession number (also called 'Entry') from column 1 of the results table, and use this to retrieve the entry into PHI-Canto on the gene entry page. **Caution:** Do not confuse the 'Entry' column with the 'Entry name' column. PHI-Canto uses the accession number to retrieve details (such as the gene name, gene product, and organism). If PHI-Canto is unable to find your entry, check for typos (e.g., 0 for O), ensure you are using the 'entry' not 'entry name', and check that your accession is from UniProtKB, not UniParc.

Figure 5 – figure supplement 1. Resources relied upon by PHI-Canto.

Category	Resource	URL	Description of use
Databases	PHI-base	http://www.phi-base.org/	To display Pathogen-Host Interaction annotation data.
	UniProtKB	https://www.uniprot.org/uniprot/	To identify the gene product under annotation.
	NCBI Taxonomy	https://www.ncbi.nlm.nih.gov/taxonomy	To identify the species of the organism being annotated.
PHI-base controlled vocabularies	PHIDO	https://raw.githubusercontent.com/PHI-base/phido/master/phido.obo	To annotate wild type metagenotypes with the 'disease caused'.
	Pathogen Host Interactions Experimental Conditions Ontology	https://raw.githubusercontent.com/PHI-base/phi-eco/master/phi-eco.obo	To annotate experimental conditions on phenotype annotations.
OBO ontologies	BioGrid	https://raw.githubusercontent.com/BioGRID/BioGRID-Ontologies/master/BioGRIDExperimentalSystems.obo	To annotate physical interactions.
	BRENDA Tissue Ontology	http://purl.obolibrary.org/obo/bto.obo	To annotate the extension 'infected tissue'.
	ChEBI ontology	http://purl.obolibrary.org/obo/chebi.obo	To annotate extensions for Gene Ontology and RNA level annotations.
	Gene Ontology	http://purl.obolibrary.org/obo/go/go-basic.obo	To annotate gene products.
	Pathogen-Host Interactions Phenotype Ontology	http://purl.obolibrary.org/obo/phipo.obo	To annotate genotypes with the observed pathogen host interaction phenotype.
	PSI-Mod Mass Modifications Ontology	http://purl.obolibrary.org/obo/mod.obo	To annotate protein chemical modifications.
	Relations Ontology	http://purl.obolibrary.org/obo/ro.obo	Contains essential relations for OBO ontologies; needed to initialize Canto.
	Sequence Ontology	http://purl.obolibrary.org/obo/so.obo	To annotate extensions for Gene Ontology annotations.