

A Coordinated Approach by Public Domain Bioinformatics Resources to Aid the Fight Against Alzheimer's Disease Through Expert Curation of Key Protein Targets

Lionel Breuza^a, Cecilia N. Arighi^{b,c}, Ghislaine Argoud-Puy^a, Cristina Casals-Casas^a, Anne Estreicher^a, Maria Livia Famiglietti^a, George Georghiou^d, Arnaud Gos^a, Nadine Gruaz-Gumowski^a, Ursula Hinz^a, Nevila Hyka-Nouspikel^a, Barbara Kramarz^c, Ruth C. Lovering^e, Yvonne Lussi^d, Michele Magrane^d, Patrick Masson^a, Livia Perfetto^d, Sylvain Poux^a, Milagros Rodriguez-Lopez^d, Christian Stoeckert^f, Shyamala Sundaram^a, Li-San Wang^f, Elizabeth Wu^g, Sandra Orchard^{d,*} and IMEx Consortium, UniProt Consortium

^aSwiss-Prot Group, SIB Swiss Institute of Bioinformatics, Centre Medical Universitaire, Geneva, Switzerland

^bProtein Information Resource, Georgetown University Medical Center, Washington, DC, USA

^cProtein Information Resource, University of Delaware, Newark, DE, USA

^dEuropean Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome Trust Campus, Hinxton, Cambridge, UK

^eFunctional Gene Annotation, Preclinical and Fundamental Science, Institute of Cardiovascular Science, University College London (UCL), London, UK

^fPerelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA

^gAlzforum, Cambridge, MA, USA

Accepted 5 June 2020

Abstract.

Background: The analysis and interpretation of data generated from patient-derived clinical samples relies on access to high-quality bioinformatics resources. These are maintained and updated by expert curators extracting knowledge from unstructured biological data described in free-text journal articles and converting this into more structured, computationally-accessible forms. This enables analyses such as functional enrichment of sets of genes/proteins using the Gene Ontology, and makes the searching of data more productive by managing issues such as gene/protein name synonyms, identifier mapping, and data quality.

Objective: To undertake a coordinated annotation update of key public-domain resources to better support Alzheimer's disease research.

Methods: We have systematically identified target proteins critical to disease process, in part by accessing informed input from the clinical research community.

Results: Data from 954 papers has been added to the UniProtKB, Gene Ontology, and the International Molecular Exchange Consortium (IMEx) databases, with 299 human proteins and 279 orthologs updated in UniProtKB. 7,45 binary interactions were added to the IMEx human molecular interaction dataset.

*Correspondence to: Sandra Orchard, European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-

EBI), Wellcome Trust Campus, Hinxton, Cambridge CB10 1SD, UK. E-mail: orchard@ebi.ac.uk. ORCID: 0000-0002-8878-3972

Conclusion: This represents a significant enhancement in the expert curated data pertinent to Alzheimer's disease available in a number of biomedical databases. Relevant protein entries have been updated in UniProtKB and concomitantly in the Gene Ontology. Molecular interaction networks have been significantly extended in the IMEx Consortium dataset and a set of reference protein complexes created. All the resources described are open-source and freely available to the research community and we provide examples of how these data could be exploited by researchers.

Keywords: Alzheimer's disease, Cytoscape network analysis, data curation, database, neurobiology, protein

INTRODUCTION

Alzheimer's disease (AD) is a progressive neurodegenerative disease characterized by loss of memory, inability to process new information, loss of language function, a disturbed perception of space, inability to do calculations, indifference, depression, delusions, and eventually death. Inheritable AD (familial AD) represents less than 5% of AD cases of which 10–15% have a family history of autosomal dominant inheritance; whereas the more common, sporadic, AD with complex polygenic risk inheritance accounts for more than 90% of cases [1]. Worldwide, at least 50 million people are currently believed to be living with AD or other dementias and this number could exceed 152 million by 2050 (<https://www.who.int/news-room/fact-sheets/detail/dementia>). The global cost of AD and dementia is estimated to be \$605 billion, which is equivalent to 1% of the entire world's gross domestic product. Globally, governments and medical charities spend millions of taxpayer and fundraiser dollars on biomedical research into this condition. It is therefore critical that the data generated by AD research is collated, organized and available in data resources and tools to increase the pace of discovery and innovation.

AD is a complex disease which needs to be studied at many levels, from the molecular mechanisms to the cellular composition and physiology of the brain [2]. Conditions such as vascular damage and neuroinflammation are also believed to play important roles in disease initiation and progression. Our current understanding of the causes, risk factors, and sub-types of these devastating conditions have been reviewed extensively elsewhere (for example, [2–4]) and they are not the subject of this manuscript. However, a number of key processes known to play a role in disease etiology and progression are briefly described to showcase the representation of selected proteins in UniProtKB and demonstrate how users can access information about both physiological and pathological aspects of the molecules.

Central to AD disease pathology are two processes: the extracellular formation of senile plaques in the grey matter of the brain which are primarily composed of amyloid- β precursor protein (APP)-derived amyloid- β ($A\beta$) [5, 6], and intracellular accumulation of hyperphosphorylated tau/Microtubule-associated protein tau (MAPT) protein to form neurofibrillary tangles [7, 8]. $A\beta$ oligomers are believed to contribute to cell death by interfering with neuron-to-neuron communication at synapses [9] and restricting the source of oxygen and nutrients [10], while tau tangles block the transport of nutrients and other essential molecules inside neurons [11]. Whilst the relationship between $A\beta$ and tau in AD is not fully understood, abnormal species of tau protein are believed to spread in a 'prion-like' manner between cells and its uptake may be potentiated by extracellular $A\beta$ [12, 13]. $A\beta$ peptides can be cleared intracellularly by microglia and other cell types [14–16], by transcytosis across the blood-brain barrier [17, 18], or by $A\beta$ degrading enzymes, such as insulin-degrading enzyme (IDE) and neprilysin (MME) [19, 20]. Tau has been shown to be degraded via the ubiquitin-proteasome system as well as the autophagy lysosome system [21]. Disorders in clearance of $A\beta$ and tau play a key role in the development of neurodegenerative disorders such as AD while overloading of the microglial system results in chronic inflammation [22, 23]. However, evidence has been emerging that aggregation of $A\beta$ and tau may not be the underlying causes of disease, but may be the outcome of perturbations in cellular homeostasis in the brain, occurring years to decades prior to disease onset [2, 24]. Normal brain function may be compromised by the decreased ability of the brain to metabolize glucose and aberrant lipid metabolism, such as sluggish cholesterol transport [25]. To date, over 350 human proteins have been associated with the development of AD as researchers move toward an understanding of the underlying cellular mechanisms that drive the formation of the protein aggregates and the downstream effect these have on the brain.

119 The analysis and interpretation of data gener-
120 ated from increasing large-scale examination of
121 patient-derived clinical samples relies on access to
122 high-quality bioinformatics resources. The scientific
123 content of these resources is maintained and updated
124 by professional biocurators who extract knowledge
125 from unstructured biological data described in free-
126 text journal articles and convert it into both more
127 easily digestible, high-level summaries and a struc-
128 tured, computable form. The latter both enables
129 large-scale data analyses, for example functional
130 enrichment of sets of genes/proteins using the Gene
131 Ontology (GO) [26, 27], and also helps to make
132 the searching of data more productive by managing
133 issues such as the problems caused by gene/protein
134 name synonyms, identifier mapping, and minimizing
135 the effect of poor quality, redundant, or mislead-
136 ing data. The work of these data resources helps
137 researchers overcome known bottlenecks in data
138 analysis, namely the time spent in discovering and
139 collating required information, manually verifying it,
140 and integrating it into analysis pipelines [28]. We here
141 describe a coordinated approach to updating key pub-
142 lic domain resources with the aim of supporting AD
143 research, starting with the update of genes/proteins
144 with a known role in AD biology. Accessing informed
145 input from the clinical research community was an
146 essential part of this process and was critical in
147 defining where curation effort was focused. We also
148 illustrate the way this coordinated update can be used
149 by researchers to answer questions pertaining to the
150 complex etiology of AD.

151 METHODS AND MATERIALS

152 *Identifying disease-related proteins*

153 A recent initiative by the UniProt Knowledgebase
154 of protein sequences and annotations [29] to update
155 the proteins which play a role in the initiation and
156 development of AD, coordinated with the curation of
157 their interactions and the complexes they form, has
158 been funded by the NIH National Institute on Aging
159 (NIA). At the start of this annotation project, cura-
160 tors were faced with two main problems—an accurate
161 description of the various forms of AD and identifi-
162 cation and prioritization of the proteins associated
163 with the disease. AD is generally classified into early
164 and late-onset forms, with genetic variants or risk
165 alleles [30] associated with each condition provid-
166 ing a further sub-classification. In order to identify
167 key AD-related proteins appropriate for update and

reannotation, UniProt curators reached out to mem- 168
bers of the AD clinical and research communities, 169
leveraging contacts made through the NIH NIA 170
programs and a collaboration with the Alzheimer's 171
Research UK (ARUK) funded GO project at Uni- 172
versity College London (UCL) [31, 32]. Workshops 173
were organized to help database providers under- 174
stand how their resources are used by the research 175
community, and conversely for the research com- 176
munity to directly input into the curation process. 177
Attendees were asked to identify proteins which 178
played a key role in the disease, or which had 179
been associated with disease even if a clear molec- 180
ular mechanism explaining this association had yet 181
to be identified. Additional candidates were pro- 182
vided by Alzforum (<http://www.alzforum.org>), the 183
Agora portal (<https://agora.ampadportal.org>), col- 184
lected from targeted research groups, and from 185
literature searching. The main pathway resource 186
consulted was WikiPathways which provided an 187
overview of the disease process (<http://www.wikipathways.org/index.php/Pathway:WP2059>). Drug target 188
resources included the ChEMBL database [33] and 189
the OpenTargets platform [34], taking only high scor- 190
ing (0.8 to 1) targets associated with AD from the 191
latter. To build the AD-centric protein-protein inter- 192
action network, data was downloaded from the IntAct 193
molecular interaction database [35], limited to inter- 194
actors with an MIscore of >0.45 (see explanation 195
below). Proteins were prioritized for curation follow- 196
ing a ranking system, i.e., 1) proteins known to play 197
a functional role in AD pathways and known drug tar- 198
gets for AD, 2) proteins known to have an association 199
to AD, e.g., through a genome wide association study 200
(GWAS) study but for which a molecular mechanism 201
has yet to be identified, and 3) proteins that physi- 202
cally interact with those defined in (1) or (2). A copy 203
of this list, as of UniProt release 2019_10 is available 204
as Supplementary Table 1. 205
206

207 *Protein annotation*

208 Data from selected publications were trans-
209 ferred into the UniProtKB, GO, IntAct molecular
210 interaction, and the Complex Portal databases, as
211 appropriate, as previously described [26, 32, 35–37].

212 *Producing an AD-centric molecular interaction 213 network*

214 Seed proteins were identified by searching the
215 UniProt website (Release 2019_08) for reviewed

entries containing the keyword 'Alzheimer disease'. (keyword: "Alzheimer disease [KW-0026]" AND reviewed: yes). As this keyword is only added to human entries, there was no need to further restrict the search by species. The final list is available in Supplementary Table 2.

Interactors of this list of proteins were obtained from IntAct using the PSICQUIC client app in Cytoscape Version: 3.7.1 [38]. To return an isoform- and post-processed chain- specific network the following query was used: (id:P37840* OR id:P49810* OR id:P49768* OR id:O14672* OR id:P03886* OR id:Q8IZY2* OR id:Q16643* OR id:P02649* OR id:P05067* OR id:Q92673* OR id:P03891* OR id:O95185*) AND annot: "imex curation".

This gave a raw network containing 1461 nodes and 2671 edges.

The network was then filtered to: a) remove non-human interactors; b) remove duplicated interactions; c) select interactions having MIscore > 0.45.

A MIscore of >0.45 can only be achieved by interacting pairs having at least a single interaction evidence showing that the two molecules directly interact or two or more evidences of a physical interaction. The filtered isoform- and post-processed chain-specific Network contains 152 nodes and 277 edges.

To enable users to access a detailed view of this network, a copy has been deposited at the NDEX data repository (<http://www.ndexbio.org/#/network/49e43d68-939b-11ea-aaef-0ac135e8bacf>) [39]. Users may alternatively download an updated set of the data used to derive an AD-focused interaction network by pasting the query annot: "dataset:Alzheimers" into the IntAct website (<http://www.ebi.ac.uk/intact>). Using the Advanced Search capabilities will enable further filtering of the results of this query.

To perform the ClueGO functional enrichment analysis, all isoform and post-processed chains were collapsed to the canonical identifiers in UniProtKB, all leaves (proteins not directly connected in the network) were removed and the complexes were demerged into protein subunits. The Cytoscape APP ClueGO version 2.5.0 [40] was then used, implementing the following parameters:

Organism analyzed: Homo Sapiens [9606].

Identifier types used: UniProtKB

#Genes in custom reference set: 3001 human proteins extracted from UniProt having tissue-specificity='brain'

Ontology used: GO_BiologicalProcess-EBI-QuickGO-GOA_20.11.2017_00h00 and REACTOME_Reactions_20.11.2017

Evidence codes used: All

Statistical Test Used=Enrichment (Right-sided hypergeometric test), Correction Method Used= Bonferroni step down

Min GO Level= 8

Max GO Level= 20

GO Fusion= true

GO Group= true

Kappa Score Threshold= 0.4

Over View Term= SmallestPValue

Group By Kappa Statistics= true

Initial Group Size= 1

Sharing Group Percentage= 60.0

RESULTS

All known human protein-coding genes have been curated by experts within the UniProtKB database (<http://www.uniprot.org>) with, as far as possible, all the protein products encoded by one gene described in a single reviewed entry [29, 36]. Each entry groups all the protein isoforms expressed by that gene, with positional features such as binding domains, post-translational modifications and amino-acid variants mapped to a representative sequence. Isoforms yet to be integrated are maintained in unreviewed entries but are accessible as part of the complete human proteome reference set (UniProt Proteome UP000005640) and can also be viewed in the corresponding reviewed entry on the website as a result of an automatic gene-centric mapping. Expert curators summarize knowledge extracted from biomedical literature in sections describing different aspects of protein biology relevant to those gene products, these can include function, enzymatic activity, subcellular location, and links to disease conditions. For example, over the period of this annotation project PSEN1 (UniProtKB P49768) had data from 43 publications added to its entry in UniProtKB, enhancing the 'Function' section, and including details of the functional roles played by specific domains within the protein. Information on disease linked variants and the effects of point mutations on protein behavior were also added.

Proteins do not operate in isolation and details of their interactions with other molecules are manually curated by the IMEx Consortium of interaction databases (<http://www.imexconsortium.org>) [41] via

315 the IntAct database [35], from where a filtered subset
316 of high confidence binary protein-protein interactions
317 is imported back into the 'Interaction' section of the
318 corresponding UniProtKB entries. Proteins also form
319 higher-order, functional assemblies and descriptions
320 of stable protein complexes are curated into the Com-
321 plex Portal (<http://www.ebi.ac.uk/complexportal>),
322 giving details of complex content, stoichiometry,
323 and topology in addition to function and 3D struc-
324 ture, when available [37]. Again, these data can be
325 accessed from the appropriate UniProtKB records.
326 In parallel, biocurators link these proteins and pro-
327 tein complexes to appropriate GO terms describing
328 their biological function, the cellular processes in
329 which they play a role, and the cellular compartment
330 in which they are found. The GO is a biomedical
331 ontology which describes these aspects of protein
332 behavior in a consistent and computer-accessible
333 manner [26, 27]. Linking gene products to GO terms
334 means that researchers can use the resulting annota-
335 tions to interpret high-throughput datasets using GO
336 term enrichment.

337 The NIA-funded annotation project resulted in data
338 from 954 papers being added to the UniProtKB, GO,
339 and IMEx databases, with 299 human proteins and
340 279 orthologues updated in UniProtKB. 7,045 binary
341 interactions were added to the IMEx human molecu-
342 lar interaction dataset.

343 *Understanding the function of AD-associated* 344 *proteins*

345 UniProt curators provide high-quality literature
346 sourced annotations for experimentally characterized
347 proteins across diverse protein families. These data
348 are presented both in free text fields and in struc-
349 tured mappings to the underlying protein sequence
350 to enable users to understand how, for example, a
351 post-translational modification to a specific residue
352 can drive a change in protein behavior. The pro-
353 teins identified by AD domain experts were subjected
354 to an intense literature review and corresponding
355 update of the relevant annotation fields in order to
356 help researchers understand both the physiological
357 role these entities play in a cell, and how this relates
358 to the pathological disease condition. As described
359 above, this includes a full review of both protein iso-
360 forms and protein chains formed by post-translational
361 processing of the full-length gene product. This is
362 particularly important in the case of AD-related
363 proteins as amyloid plaque formation is a conse-
364 quence of dysregulated protein cleavage [42]. APP

(UniProtKB P05067) is a ubiquitously expressed type
I transmembrane protein which functions as a cell sur-
face receptor with roles in neurite growth, neuronal
adhesion, and axonogenesis. The protein consists of
a large ectodomain, a single membrane spanning
domain and a short cytoplasmic tail. The ectodomain
comprises two highly conserved E1 and E2 domains,
involved in metal (copper and zinc) and heparin
binding. APP undergoes extensive post-translational
modification and proteolytic processing to generate
peptide fragments. The cleavage products of APP are
all described at the residue level in the UniProtKB
database, with stable identifiers allowing unambigu-
ous recognition of each proteoform when described
(Fig. 1).

As detailed in the appropriate UniProtKB
records, APP processing is initiated either by α -
secretase/ADAM10 (UniProtKB O14672) cleavage
within the A β region, or by β -secretase (BACE1/2,
UniProtKB P56817/Q9Y5Z0) cleavage at the N-
terminus of A β , leading to the secretion of large
soluble ectodomains, termed soluble APP α (APP α ,
UniProtKB PRO_0000000089) and soluble APP β
(APP β , UniProtKB PRO_0000000090), respec-
tively. Subsequent processing of the C-terminal
fragments by the γ -secretase complex (Com-
plex Portal:CPX-2176/CPX-4231/CPX-4232/CPX-
4233), as well as processing along non-canonical
pathways, results in numerous fragments, which have
different and partially opposite functional properties.
During amyloidogenic processing, APP is sequen-
tially cleaved by β - and γ -secretases to mainly
generate A β_{40} (UniProtKB PRO_0000000093), and
A β_{42} (UniProtKB PRO_0000000092) fragments.

Many of the AD-associated proteins prioritized for
update (Supplementary Table 1) are enzymes, which
may be responsible for the proteolytic processing
of longer protein chains as described above, cataly-
sis of metabolic reactions, or generation/removal
of post-translational modification sites. Enzymatic
function is now described in UniProtKB using
Rhea (<http://www.rhea-db.org>), a comprehensive and
non-redundant resource of expert-curated biochem-
ical reactions [43], as a vocabulary to annotate
and represent enzyme-catalyzed reactions. Rhea
uses the ChEBI (Chemical Entities of Biological
Interest) ontology to describe reaction participants,
their chemical structures, and chemical transforma-
tions [44]. Additional small molecule interactions,
such as cofactor binding sites are also described
within UniProt using ChEBI. Sophisticated searches
within UniProtKB now allow the researcher to

PTM / Processing¹

Molecule processing					
Feature key	Position(s)	Description	Actions	Graphical view	Length
Signal peptide ¹	1 – 17	3 Publications	Add BLAST		17
Chain ¹ (PRO_0000000088)	18 – 770	Amyloid-beta precursor protein	Add BLAST		753
Chain ¹ (PRO_0000000089)	18 – 687	Soluble APP-alpha	Add BLAST		670
Chain ¹ (PRO_0000000090)	18 – 671	Soluble APP-beta	Add BLAST		654
Chain ¹ (PRO_0000381966)	18 – 286	N-APP	Add BLAST		269
Chain ¹ (PRO_0000000091)	672 – 770	C99	Add BLAST		99
Chain ¹ (PRO_0000000092)	672 – 713	Amyloid-beta protein 42	Add BLAST		42
Chain ¹ (PRO_0000000093)	672 – 711	Amyloid-beta protein 40	Add BLAST		40
Chain ¹ (PRO_0000000094)	688 – 770	C83	Add BLAST		83
Peptide ¹ (PRO_0000000096)	688 – 713	P3(42)	Add BLAST		26
Peptide ¹ (PRO_0000000096)	688 – 711	P3(40)	Add BLAST		24
Chain ¹ (PRO_0000384574)	691 – 770	C80	Add BLAST		80
Chain ¹ (PRO_0000000097)	712 – 770	Gamma-secretase C-terminal fragment 59	Add BLAST		59
Chain ¹ (PRO_0000000098)	714 – 770	Gamma-secretase C-terminal fragment 57	Add BLAST		57
Chain ¹ (PRO_0000000099)	721 – 770	Gamma-secretase C-terminal fragment 50	Add BLAST		50
Chain ¹ (PRO_0000000100)	740 – 770	C31	Add BLAST		31

Fig. 1. Screenshot showing the UniProtKB description of the products of amyloid-beta precursor protein post-transcriptional modifications and processing. This information is available in the UniProtKB P05067 entry for amyloid-beta precursor protein (APP).

417 identify metabolic networks and predict new path-
 418 ways for drug production. For example, alterations
 419 in sphingolipid metabolism have been detected in
 420 AD, with levels of SPHK1 (UniProtKB Q9NYA1)
 421 downregulated and, conversely, levels of SPHK2
 422 (UniProtKBQ9NRA0) upregulated [45]. Both entries
 423 for these proteins have been updated in UniProtKB,
 424 where it is now possible to visualize the chemical
 425 reaction, balanced for mass and charge (at an arbi-
 426 trary pH of 7.3) as described by Rhea, and cofactors
 427 linked to the corresponding entry in ChEBI (Fig. 2).

428 Tau/MAPT (UniProtKB P10636) is a microtubule-
 429 associated protein predominantly expressed in the
 430 axons of neurons [46]. Tau is a naturally unfolded
 431 protein with an extended structure; however, in AD
 432 brains, tau is accumulated in a hyperphosphory-
 433 lated state in a unique filamentous structure, paired
 434 helical filaments of 10 nm diameter with 80 nm peri-
 435 odicity [47]. The phosphorylation of tau regulates
 436 both its functional ability to assemble and stabilize
 437 microtubules and also its pathological structure [48],
 438 and the 441 amino acid isoform of tau (UniProt-
 439 tKB P10636-8) has 45 serine, 35 threonine, and 5
 440 tyrosine residues, resulting in a total of 85 poten-
 441 tial phosphorylation sites [49]. CDK5 (UniProtKB
 442 Q00535) is one enzyme known to play a role in
 443 the phosphorylation of tau [50], priming tau for
 444 further phosphorylation events by the hierarchical
 445 kinase GSK3B (UniProtKB P49841) by modify-
 446 ing an upstream +4 (or +3) site, (S/T)xx(x)p(S/T).
 447 Again, this chemical reaction has been updated in

448 UniProt (Fig. 2B), where it is also possible to
 449 identify the resulting phosphorylated residues in
 450 the corresponding entry for tau. CDK5 is activated
 451 by p35/CDK5R1 (Q15078), the resulting complex
 452 (Complex Portal:CPX-2201) then being recruited to
 453 membranes via the N-terminal p35 myristoylation
 454 site (Fig. 2C) [51]. p35/CDK5R1 is a protein with a
 455 short-life span which is cleaved by calpain (Complex
 456 Portal:CPX-2674/CPX-4302) into a p25 C-terminal
 457 fragment (UniProtKB PRO_0000004795) when neu-
 458 rons suffer from stress or encounter death signals.
 459 p25/CDK5R1 has a longer half-life and this com-
 460 plex (Complex Portal:CPX-3142) dissociates from
 461 the plasma membrane into the nucleus, where it can
 462 phosphorylate additional proteins [52].

463 Linking amino acid variation to functional 464 consequence

465 AD-causing mutations in APP (UniProtKB
 466 P05067), PSEN1 (UniProtKB P49768), and PSEN2
 467 (UniProtKB P49810) affect the generation of A β
 468 peptides, changing the relative ratio of A β ₄₂ to
 469 A β ₄₀ peptide [53]. The longer A β ₄₂ peptides seem
 470 to be more prone to aggregation, and increased
 471 ratios of A β ₄₂/A β ₄₀ are thought to play a role in
 472 AD pathogenesis. It is therefore important to doc-
 473 ument all APP, PSEN1, and PSEN2 variants that
 474 lead to a change in this ratio. About 1% of AD
 475 cases develop as a result of mutations within APP or
 476 the genes encoding the PSEN1 and PSEN2 proteins

Entry	Entry name	Protein names	Gene names	Organism	Length	Catalytic activity	Cofactor
Q9NYA1	SPHK1_HUMAN	Sphingosine kinase 1	SPHK1 SK1, SPHK, SPK	Homo sapiens (Human)	384	<ul style="list-style-type: none"> a sphingoid base + ATP = a sphingoid 1-phosphate + ADP + H⁺ 5 Publications EC:2.7.1.91 5 Publications Source: Rhea. acetyl-CoA + L-seryl-[protein] = CoA + O-acetyl-L-seryl-[protein] By similarity Source: Rhea. ATP + sphinganine = ADP + H⁺ + sphinganine 1-phosphate 2 Publications EC:2.7.1.91 2 Publications Source: Rhea. ATP + sphing-4-ene = ADP + H⁺ + sphing-4-ene 1-phosphate 2 Publications EC:2.7.1.91 2 Publications Source: Rhea. 1-O-hexadecyl-2-amino-sn-glycerol + ATP = 1-O-hexadecyl-2-desoxy-2-amino-sn-glycerol-3-phosphate + ADP + H⁺ 1 Publication Source: Rhea. 	Mg ²⁺ 1 Publication

Fig. 2. Screenshot showing the results of a UniProtKB search for human Sphingosine kinase 1. The UniProtKB database was queried for the term 'SPHK1'. The top hit (human) in the results table is displayed. It is possible to customize this view to select additional data fields from the UniProt record, in this case the column options 'cofactor' and 'catalytic activity' were added to the results table.

477 present in the γ -secretase complex; however, those
 478 inheriting a known AD-associated APP or PSEN1
 479 variant will develop the disease, whereas a slightly
 480 lower risk (95%) is associated with inheriting a
 481 known AD variant in PSEN2 [54]. Individuals with
 482 AD mutations in any of these three genes tend to
 483 develop early-onset disease, with symptoms devel-
 484 oping before the age of 65, sometimes as early
 485 as age 30. Understanding how a genetic varia-
 486 tion changes protein function or expression levels
 487 is essential for our understanding of genetic dis-
 488 ease and the ability to identify those variants which
 489 are causal. UniProtKB curators capture nonsynony-
 490 mous variants described in the literature with, when
 491 available, detail on the phenotypic or pathogenic
 492 consequences on the amino acid change. UniProt
 493 also receives input (publications and suggested anno-
 494 tations) from expert groups, e.g., Alzforum, who
 495 collects detailed variant information about AD pro-
 496 teins from the literature. To date, UniProtKB records
 497 contain information on over 30,000 variants linked
 498 to Mendelian diseases in more than 13,000 human
 499 protein sequence records [55] and work is ongo-
 500 ing to standardize variant interpretations through
 501 the incorporation of American College of Medical
 502 Genetics and Genomics (ACMG) guidelines and the
 503 ClinGen pathogenicity calculator into the curation

workflow. Cross references to variant resources
 such as dbSNP (<http://www.ncbi.nlm.nih.gov/snp/>)
 and Ensembl (<http://www.ensembl.org/>), and disease-
 specific databases such as NIAGADs (<https://www.niagads.org/>) are added. Additional variant data is
 imported from large-scale studies such as 1000
 Genomes and ExAC, and again mapped to the pro-
 tein sequence and made available via the Proteins API
 (<http://www.ebi.ac.uk/proteins/api/doc/>).

UniProtKB acts as an integrative layer, enabling
 users to align genomic variants with enzyme active
 sites, modified residues, the phenotypic consequence
 of site-directed mutagenesis and binding domains
 mapped to the residue level. An exact mapping of
 the Ensembl translation to a UniProtKB sequence
 enables the calculation of UniProtKB positional
 annotations to their genomic coordinates and these
 mappings are continually reviewed and updated by
 both UniProt and Ensembl curation teams [56].
 Thirty-four different positional annotation types are
 currently aligned with the genome sequence. An addi-
 tional 17,371 mutations which map to the genome
 have been supplied by the IMEx Consortium which
 captures the effects of point mutations on molecu-
 lar interactions, using controlled vocabulary terms to
 describe whether these increase, disrupt, or cause an
 interaction to occur [57]. Again, these site-directed

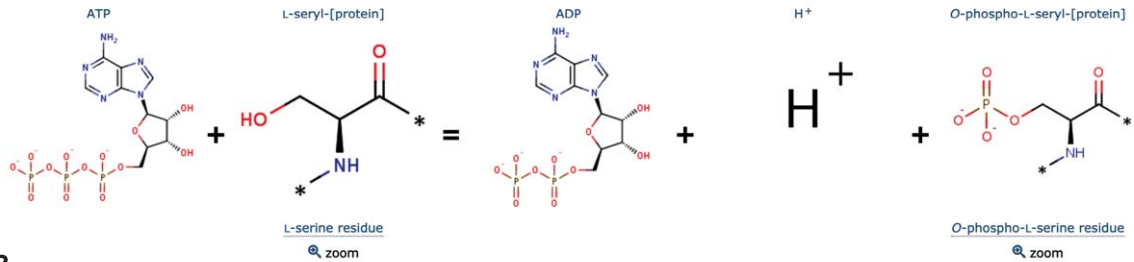
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530

A

Protein | **Cyclin-dependent-like kinase 5**
 Gene | **CDK5**
 Organism | *Homo sapiens (Human)*
 Status | Reviewed - Annotation score: - Experimental evidence at protein level¹

Catalytic activity¹

- ATP + L-seryl-[protein] = ADP + H⁺ + O-phospho-L-seryl-[protein]
 EC:2.7.11.1
 Source: Rhea. [← Hide](#)



B

Cyclin-dependent protein kinase 5 holoenzyme complex, p35 variant

ComplexAc: CPX-2201

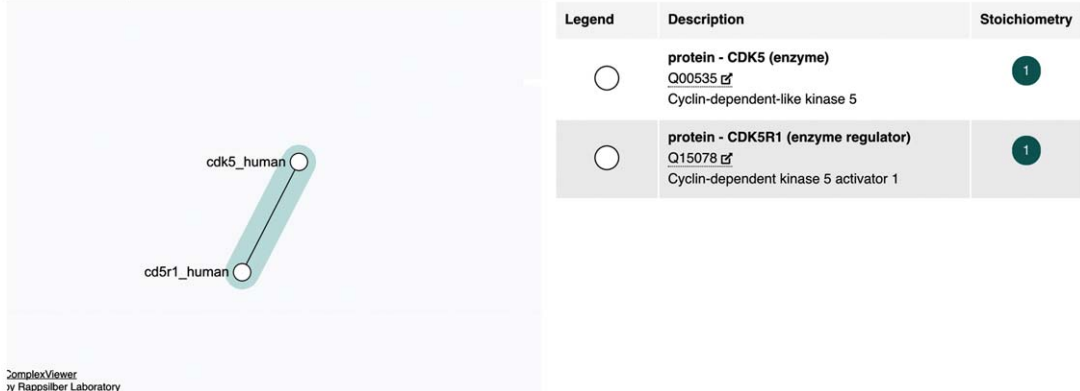
Homo sapiens; 9606

Fig. 3. Representation of CDK5 in UniProtKB and the Complex Portal. A) The representation of CDK5 catalytic activity by Rhea within the UniProtKB entry. B) The Complex Portal display of the CDK5-p35/CDK5R1 complex which can be found by searching for either of the proteins, or by complex name.

531 mutations have been mapped to the underlying
 532 UniProt protein sequence and can be used to understand
 533 the effect a genomic variant may have on a
 534 local protein interaction network. Further to this, in
 535 collaboration with PDBe through the Structure Integration
 536 with Function, Taxonomy, and Sequences resource (SIFTS;
 537 <http://pdbe.org/sifts/>), UniProtKB maps
 538 between protein structure and protein sequence, so that a
 539 knowledge of protein conformation can

contribute to an understanding of protein function
 [58]. These data are all displayed in UniProtKB using the
 ProtVista visualization tool [59] which allows the graphical
 alignment of sequence feature data to the linear protein
 sequence and from there to the 3D structure (Fig. 4).

Late-onset AD is observed in >90% of patients, and the
 APOE (UniProtKB P02649) allele E4 is strongly associated
 with these cases. APOE is a plasma

540
 541
 542
 543
 544
 545
 546
 547
 548
 549

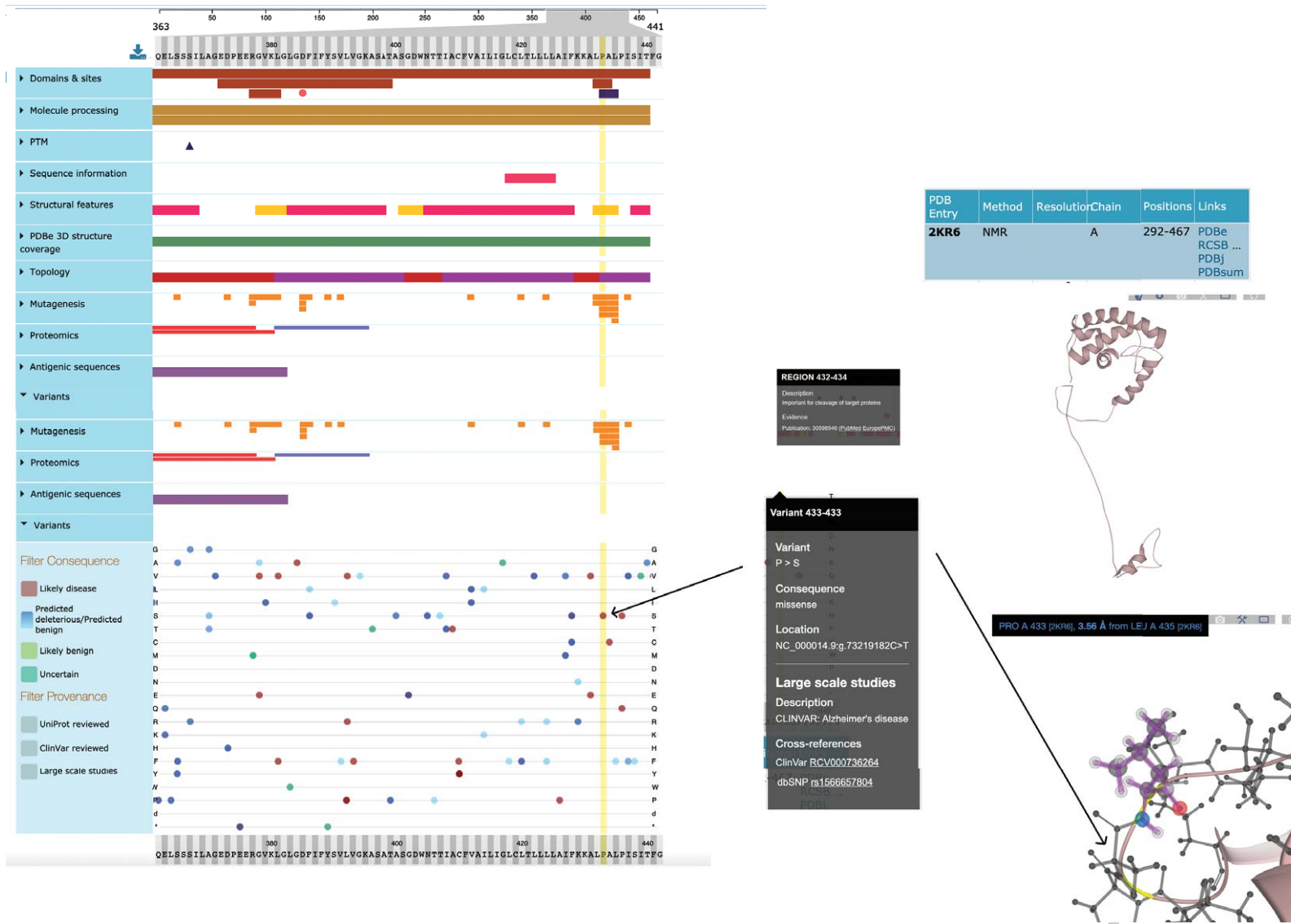


Fig. 4. Simplified view of ProtVista for Human PSEN1 (UniProtKB P49768). To investigate the effect of a specific variant (p.Pro433Ser) of human PSEN1 protein, the user can look at its potential effect on active sites and domain. Clicking on the variant at position 433 shows it to be positioned in the PAL domain, required for normal active site conformation and also in a region important for cleavage of this protein. The position of this variant is also highlighted in the NMR structure of this protein.

lipoprotein which transports lipids between cells and tissues. Abnormal cholesterol metabolism associated with allele E4 is believed to mediate cell type-specific AD pathology, including A β upregulation and impaired synaptic function in neurons, reduced synapse elimination activity in astrocytes, impaired remyelination in oligodendrocytes, and A β accumulation and inflammatory response in microglia [60]. The most common allele in the human population, and that present on the reference genome GRCh38, APOE*3 is the displayed sequence in the UniProtKB entry, with all three possible alleles fully described in the Polymorphism section of the entry. Sequence variants, single amino acid polymorphisms, and other sequence annotations, have then been described relative to that allele with the alignment of the APOE sequence to the reference genome then allowing the integration of genomic and protein data. In the recent curation project, APOE had information from 40 new references added to the entry.

Enabling functional Insights Into large-scale AD datasets through network analysis

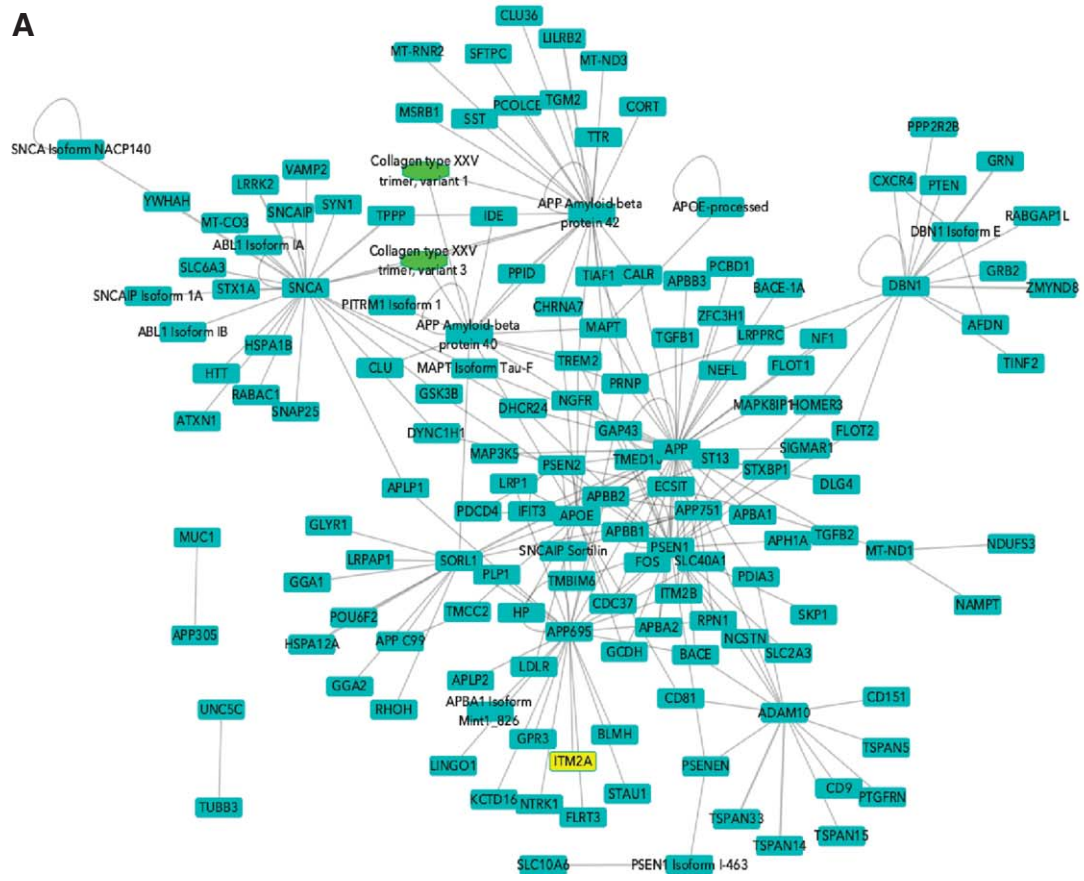
AD is not a single disease but a number of separately-triggered conditions [61] which share the same pathological phenotype, suggesting that these conditions may have many downstream processes in common. Understanding how proteins associated with AD are linked in the interacting network of molecules that drive cellular processes may help to identify proteins which are critical for initiating or driving the disease condition as potential therapeutic targets. Network-based analysis is a powerful technique for extracting biological insights from large datasets, enabling researchers to identify clusters of interacting molecules which participate in the same biological process or are members of the same physical complex. Protein interaction networks can help researchers understand the interconnectivity of both intra- and extracellular signaling, while studying network topology can give information about biological function and properties of the component molecules. Merging external 'omics data, such as transcriptomics, proteomics, and genome-wide association (GWA) studies, with the network can indicate tightly associated nodes of co-regulated proteins. An understanding of the processes associated with these networks can be further investigated by using GO annotations or Complex Portal data.

The IMEx Consortium curates to a detailed curation model, i.e., all aspects of an interaction

experiment, including host organism, interaction detection, and participant identification methodologies and full details of the constructs, including binding domains and the effects of site-directed mutations, are captured [41, 57]. All this information is accurately mapped to controlled vocabulary terms, in particular those described by the HUPO PSI-MI CV. Interactions are not limited to protein:protein but increasingly also include protein-small molecule, protein-protein complex, protein-ncRNA, and protein-gene interactions using identifiers from ChEBI, Complex Portal, RNA-Central (<http://www.rnacentral.org>), and Ensembl, respectively, to identify the respective entities. This enables the IMEx databases to fully capture the differences in interacting molecules observed with different APP isoforms (UniProtKB P05067-4/P05067-8 IntAct:EBI-21132406/EBI-21132308) [62] or by monomeric (UniProtKB PRO_0000000092) versus oligomeric (Complex Portal CPX-1134) A β ₄₂ (IntAct:EBI-20818781/EBI-20821761) [63]. The effects of mutagens, site directed to mimic known variants can also be described, for example the interactome of MAPT/Tau (UniProtKB P10636) p.Pro618Leu variant (dbSNP:rs63751273) with a known link to frontotemporal dementia [64], which reduces the ability of MAPT/Tau to promote microtubule assembly and accelerates aggregation of tau into filaments has been compared to that of the wild-type protein (IntAct:EBI-20800792/EBI-20799058) [65]. Data on the effect of site-directed mutations on molecular interactions is available as a downloadable file from the IntAct website (<ftp://ftp.ebi.ac.uk/pub/databases/intact/current/various/mutations.tsv>) and is also exported to the UniProtKB ProtVista viewer to provide additional understanding of how a particular amino acid variant may affect protein function.

Any protein interaction network built using current data will at best be partial, as we are far from having achieved full coverage of the human interactome. However, a more immediate concern is the quality of the networks being used for analysis, which are currently often created by combining data from many resources with little attention to the source(s) of the binary interactions and the methodology by which they were generated. The detailed curation model of the IMEx curation enables data filtering on many levels and thus enables the building of high-quality networks. The addition of AD-relevant protein interactions as a part of the curation marathon described above has enriched the interactome of AD-related proteins by several thousand binary interactions and

A



B

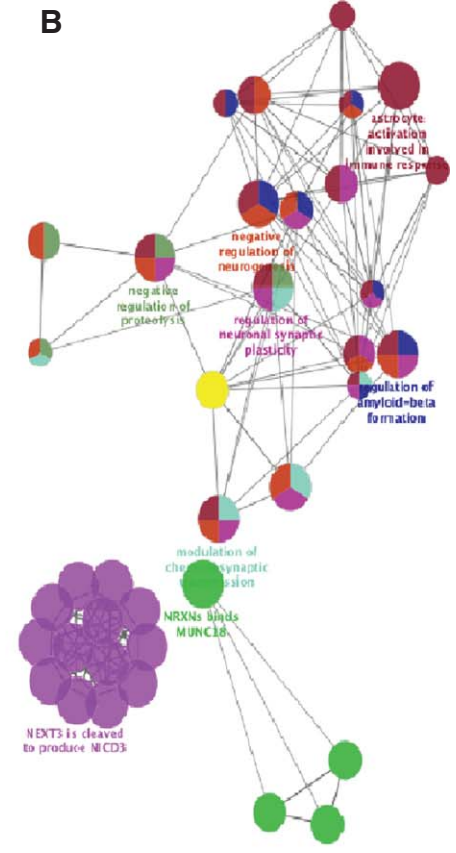


Fig. 5. Networks built from Alzheimer's disease relevant proteins. A) High-confidence network with isoforms and post-processed chains acting as distinct nodes. Seed proteins are those to which the Alzheimer Disease keyword has been added in UniProtKB. Blue squares represent proteins, green ovals represent protein complexes. Nodes have been collapsed to canonical sequence/gene level. B) ClueGO functional enrichment analysis of network shown in A.

651 is a significant addition to previous work by the
652 IMEx curators in building the APP interactome [66].
653 To demonstrate the utility of these data for AD
654 researchers, high confidence interaction networks
655 could be built using both protein interactors described
656 at the isoform/post-processed chain level and also
657 following the collapse of this level of detail to the
658 consensus sequence selected (Fig. 5A). In both cases,
659 the seed proteins were those to which the Alzheimer
660 Disease keyword has been added in UniProtKB. The
661 raw network contained 1,461 nodes and 2,671 edges.
662 This was then filtered by MI score > 0.45 [67] to
663 produce a high-confidence sub-network, restricted
664 to human-only interactions and redundant interac-
665 tion evidences were merged, reducing this to 152
666 nodes and 277 edges. Remapping isoforms and post-
667 processed chains to the canonical sequence level
668 further reduced this to 136 nodes and 179 edges.
669 This final network was analyzed using ClueGO, a
670 Cytoscape App that visualizes non-redundant bio-
671 logical GO terms for large clusters of genes in
672 a functionally grouped network (Fig. 5B). In this
673 case, the network was filtered for 'Biological Pro-
674 cess' term enrichment. Terms such as 'regulation
675 of amyloid-beta formation', 'regulation of synap-
676 tic plasticity', 'astrocyte activation' (linked to AD
677 pathology [68]), and child terms of Notch1 signaling
678 (known to be altered in AD [69]) were overexpressed
679 in comparison to a full list of human brain proteins,
680 suggesting that this is a biologically relevant network.
681 This network, and subsequent ongoing expansions to
682 the dataset, is now freely available to the research
683 community to enable network analysis of generated
684 data and can easily be extended to encompass, for
685 example, all the proteins known to be expressed in
686 the human brain by performing the relevant queries
687 on the IntAct website. This resource will facilitate
688 interrogation of large-scale GWAS, transcriptome
689 and proteomics clinical datasets, and allow users to
690 explore novel biology and enhance our understanding
691 of the disease process [70].

692 The Reactome database of curated biological path-
693 ways provides a tool for visualizing user-supplied
694 expression data as an overlay on manually curated
695 pathway diagrams [71]. Pathways are authored by
696 biologists who are recruited for their expertise in
697 the area, in this case biocurators involved with the
698 curation of AD-associated papers in UniProtKB. As
699 a result of this curation marathon, a number of
700 AD-related pathways are in the process of being cre-
701 ated and will be available to researchers as another
702 tool enabling large-scale 'omics analysis. Reactome

703 pathways can be further extended by adding IMEx
704 quality filtered protein interactions to extend out
705 the networks and these additional molecules can
706 be included in subsequent representation analysis,
707 a statistical (hypergeometric distribution) test that
708 determines whether certain Reactome pathways are
709 over-represented (enriched) in any submitted dataset.

710 *Further enhancing the Gene Ontology to improve* 711 *interpretation of AD data*

712 UniProtKB biocurators are the single largest con-
713 tributing group to GO manual annotations, both
714 as a whole but in particular for the annotation of
715 human proteins. The recent focus on AD proteins
716 has added to work by the UCL Functional Gene
717 Annotation group, funded by ARUK, to associate GO
718 terms to proteins, protein complexes, and microR-
719 NAs relevant to processes involving amyloid-beta
720 and tau, concomitantly creating many new GO terms
721 in the process to further enrich those branches of
722 the ontology relevant to neuronal biology. As a
723 proof-of-concept of the benefit of a focused anno-
724 tation effort, a functional analysis was performed by
725 Kramarz et al. in November 2018 [31] on a hippo-
726 campal proteomic dataset, identifying proteins that were
727 differentially expressed in AD versus age-matched
728 controls. Analyzing the data against the GO in 2018
729 versus an earlier version archived in 2016 showed
730 an almost doubling of enriched GO terms and high-
731 lighted new processes with a potential role in AD, for
732 example 23% of dysregulated hippocampal proteins
733 now showed a contribution to a heightened immune
734 response. The work on curating proteins and protein
735 complexes to GO terms is being continued by the
736 UniProt, Complex Portal, and UCL annotation teams,
737 while the UCL team are additionally contributing GO
738 annotation of microRNAs regulating the expression
739 of microglial AD relevant proteins [32].

740 One advantage of the improved GO representa-
741 tion of processes related to AD, is that it can be
742 used as a tool to search for lists of proteins rele-
743 vant to a particular aspect of the disease. It is now
744 widely acknowledged that neuroinflammation plays
745 a key role in the pathogenesis of AD, for example
746 through the elevation of amyloidogenesis. The list
747 of proteins involved in any inflammatory response
748 is long, but searching the UniProt or QuickGO
749 (<http://www.ebi.ac.uk/QuickGO>) websites for pro-
750 teins annotated to the GO term "neuroinflammatory
751 response" (GO:0150076) and limiting the search
752 to human proteins, retrieves a list of 42 reviewed

753 protein entries (GOA release 2020-04-22), which
754 may be connected to the disease process. The pro-
755 tein list can be downloaded from the QuickGO
756 website in CSV format, along with all the GO anno-
757 tations and publications from which the evidence was
758 extracted.

759 DISCUSSION

760 AD is a progressive brain disorder that damages
761 and destroys brain cells, leading to loss of mem-
762 ory, dysregulated brain function, and eventually death.
763 In addition to the profound human suffering caused
764 by the condition, AD and other dementias are cre-
765 ating an enormous pressure on both health care
766 systems and national budgets. To understand the
767 molecular mechanisms both triggering and subse-
768 quently driving the development of AD, researchers
769 have designed numerous high-throughput transcrip-
770 tomic, proteomic, metabolomic, and GWA studies
771 generating vast amounts of data. The subsequent
772 analyses and interpretation of the results from such
773 experiments is completely dependent on functional
774 annotation data provided by bioinformatic resources.
775 Resources such as UniProt, the GO, and the IMEx
776 molecular interaction networks enable researchers to
777 take lists of genes/proteins identified in large-scale
778 'Omics experiments and, for example, find clusters
779 of co-regulated genes which may represent processes
780 or protein complex members involved in a particular
781 process or pathway.

782 The content of these core data resources is depen-
783 dent on the work of skilled biocurators, reading and
784 evaluating the scientific literature and transferring
785 key facts to the appropriate entries. Expert manual
786 curation is undeniably expensive, but is essential
787 to make this information readily available to the
788 researcher, the clinician, and to the computational
789 biologist. By working collaboratively, contributing
790 data to multiple specialist resources and working
791 together to develop shared curation tools [72], the
792 biocuration community is taking a lead in giving fun-
793 ders the best possible return on their investment [28].
794 The AD focused biocuration project described here
795 has benefitted from governmental funding, charita-
796 ble funding, from pharmaceutical company funding
797 through a public-private partnership [32, 34] and
798 also from previously funded work into other neu-
799 rological conditions [73, 74]. While in this case
800 the shared funding pool was serendipitous, it sug-
801 gests that actively managed collaborations between

802 funding bodies could be at least equally successful
803 in increasing both the quantity and quality of infor-
804 mation freely available in biomedical databases. As
805 a result of these efforts, researchers can now access
806 299 disease-relevant human protein records updated
807 in UniProtKB (as of release 2019_10), with experi-
808 mental GO annotation also added, where possible. An
809 additional 7045 binary molecular interactions have
810 been added to the IMEx dataset, significantly increas-
811 ing the abilities of researchers to perform network
812 analysis on large-scale datasets.

813 Once the data is in these resources, it is also the
814 responsibility of database managers to ensure that
815 users can find and access it as easily as possible. The
816 UniProt Consortium is already working to release
817 a disease-specific entry point to those proteins of
818 interest which will enable researchers to navigate the
819 network of molecules that play a role in this condition
820 and easily find information on the function of each.
821 An AD portal will be the first of these released. The
822 data is also being made available through other pub-
823 lic domain biomedical resources such as the Open
824 Targets platform (<http://www.opentargets.org>) [34]
825 which integrates evidence from genetics, genomics,
826 transcriptomics, drugs, animal models, and scientific
827 literature to score and rank target-disease associations
828 for drug target identification. The UniProt Consor-
829 tium is also looking to improve the ability of both
830 scientists and clinicians to navigate from genomic
831 disease variant to amino acid polymorphism to effect
832 of protein structure and/or function with both graphi-
833 cal visualization and computational access readily
834 available. Variant data will become more structured,
835 thus making it more computationally accessible [55].
836 The value of metabolomics data derived from AD-
837 patients will be significantly enhanced by the work
838 on enhancing the content of Rhea and ChEBI, and
839 ensuring that appropriate data are incorporated into
840 UniProt and improved and updated protein sequences
841 will increase the number of identifications made by
842 mass spectrometry-based proteomics groups.

843 In conclusion, the work described above represents
844 a significant increase in the content of a number of
845 public domain resources specifically focused on the
846 molecules which play a key role in AD. Many of these
847 proteins also play a role in other neurological disor-
848 ders and are, of course, of fundamental importance
849 to the normal physiology of the brain. These ongoing
850 and future data updates will help clinical researchers
851 to provide insights into the molecular mechanisms
852 underlying the development of dementia and enable
853 more in-depth analysis of 'Omics'-level datasets, thus

854 supporting the development of novel treatments and
855 tools for early diagnosis.

856 DATA AVAILABILITY

857 UniProtKB records in which disease is caused
858 by mutations affecting the gene represented in that
859 entry can be found by searching www.uniprot.org
860 with the term “keyword:”Alzheimer disease [KW-
861 0026]”. An introduction to the QuickGO Gene
862 Ontology browser can be found at [www.ebi.ac.uk/
863 training/online/course/goa-and-quickgo-quick-tour](http://www.ebi.ac.uk/training/online/course/goa-and-quickgo-quick-tour).
864 Tutorials on how to search UniProt and use the tools
865 made available by this resource and how to access
866 data pertaining to AD in the GO are available [75,
867 76]. Data required to create AD-focused molecular
868 interaction network can be obtained by pasting the
869 query annot: “dataset:Alzheimers” into the IntAct
870 website (www.ebi.ac.uk/intact) with further details
871 on how to use this resource available at [www.ebi.ac.
872 uk/training/online/course/intact-molecular-interactio
873 ns-ebi](http://www.ebi.ac.uk/training/online/course/intact-molecular-interactions-ebi). Extensive tutorial materials on Cytoscape
874 network building and analysis are available at [https://
875 github.com/cytoscape/cytoscape-tutorials/wiki](https://github.com/cytoscape/cytoscape-tutorials/wiki), the
876 use of ClueGO is specifically described by Bindea
877 et al. [40, 77]. How to use the Complex Portal is
878 described by Meldal et al. [78].

879 ACKNOWLEDGMENTS

880 This work was supported by the National Eye
881 Institute (NEI), National Human Genome Research
882 Institute (NHGRI), National Heart, Lung, and Blood
883 Institute (NHLBI), National Institute on Aging
884 (NIA), National Institute of Allergy and Infectious
885 Diseases (NIAID), National Institute of Diabetes
886 and Digestive and Kidney Diseases (NIDDK),
887 National Institute of General Medical Sciences
888 (NIGMS), National Cancer Institute (NCI) and
889 National Institute of Mental Health (NIMH) of the
890 National Institutes of Health under Award Number
891 [U24HG007822]. Research reported in this publi-
892 cation was additionally supported by the National
893 Human Genome Research Institute (NHGRI) and
894 the National Institute on Aging (NIA) of the
895 National Institutes of Health under Award Number
896 [3U24HG007822-05S1] (the content is solely the
897 responsibility of the authors and does not necessarily
898 represent the official views of the National Institutes
899 of Health).

IntAct, the Complex Portal and other EMBL-EBI-
900 based authors also received funding from EMBL
901 core funding, Open Targets (grant agreements OTAR-
902 044 and OTAR02-048) and the Wellcome Trust
903 grant INVAR (grant ref: 212925/Z/18/Z). Authors
904 based in the Swiss-Prot Group, SIB Swiss Insti-
905 tute of Bioinformatics also receive funding from the
906 Swiss Federal Government through the State Secre-
907 tariat for Education, Research and Innovation (SERI).
908 The University College London functional anno-
909 tation team is supported by ARUK-NSG2016-13,
910 ARUK-NAS2017A-1 and the National Institute for
911 Health Research University College London Hospi-
912 tals Biomedical Research Centre

914 The authors would like to thank Dr Rina Ban-
915 dhopadhyay and Profs John Hardy and (UCL,
916 UK), Profs Nigel Hooper and David Brough
917 (U.Manchester, UK), Profs Casey Brown, Li-San
918 Wang, Christian Stoeckert (U. Penn, US), Prof.
919 Michael MacCoss (U. Washington, US), Prof. Hans-
920 Ulrich Klein (Columbia U., US), Prof. Christopher
921 Martens (U. Delaware, US), Prof. Thomas Wingo
922 (Emory U. US), Dr. Christopher Khalid-Janney
923 (Delaware State U. US) amongst others for their help
924 in identifying gene candidates for annotation.

925 Authors' disclosures available online ([https://
926 www.j-alz.com/manuscript-disclosures/20-0206r1](https://www.j-alz.com/manuscript-disclosures/20-0206r1)).

SUPPLEMENTARY MATERIAL 927

928 The supplementary material is available in the
929 electronic version of this article: [https://dx.doi.org/
930 10.3233/JAD-200206](https://dx.doi.org/10.3233/JAD-200206).

REFERENCES 931

- 932 [1] Armstrong RA (2013) What causes Alzheimer's disease?
933 *Folia Neuropathol* **51**, 169-188.
- 934 [2] De Strooper B, Karran E (2016) The cellular phase of
935 Alzheimer's disease. *Cell* **164**, 603-615.
- 936 [3] Chávez-Gutiérrez L, Szaruga M (2020) Mechanisms of neu-
937 rodegeneration - insights from familial Alzheimer's disease.
938 *Semin Cell Dev Biol*, doi: 10.1016/j.semcdb.2020.03.005
- 939 [4] Sengoku R (2020) Aging and Alzheimer's disease pathol-
940 ogy. *Neuropathology* **40**, 22-29.
- 941 [5] Hardy JA, Higgins GA (1992) Alzheimer's disease: The
942 amyloid cascade hypothesis. *Science* **256**, 184-5.
- 943 [6] Selkoe DJ, Hardy J (2016) The amyloid hypothesis of
944 Alzheimer's disease at 25 years. *EMBO Mol Med* **8**, 595-
945 608.
- 946 [7] Delacourte A, Dufosse A (1986) Alzheimer's disease: Tau
947 proteins, the promoting factors of microtubule assembly, are
948 major components of paired helical filaments. *J Neurol Sci*
949 **176**, 173-186.

- 950 [8] Nolan A, De Paula Franca Resende E, Petersen C, Neylan K, 1014
 951 Spina S, Huang E, Seeley W, Miller Z, Grinberg LT (2019) 1015
 952 Astrocytic tau deposition is frequent in typical and atypical 1016
 953 Alzheimer disease presentations. *J Neuropathol Exp Neurol* 1017
 954 **78**, 1112-1123. 1018
- 955 [9] Jarosz-Griffiths HH, Noble E, Rushworth JV, Hooper NM 1019
 956 (2016) Amyloid- β receptors: The good, the bad, and the 1020
 957 prion protein. *J Biol Chem* **291**, 3174-3183. 1021
- 958 [10] Nortley R, Korte N, Izquierdo P, Hirunpattarasilp C, 1022
 959 Mishra A, Jaunmuktane Z, Kyrargyri V, Pfeiffer T, Khen- 1023
 960 nouf L, Madry C, Gong H, Richard-Loendt A, Huang 1024
 961 W, Saito T, Saido TC, Brandner S, Sethi H, Attwell D 1025
 962 (2019) Amyloid β oligomers constrict human capillaries 1026
 963 in Alzheimer's disease via signaling to pericytes. *Science* 1027
 964 **365**, eaav9518. 1028
- 965 [11] Brion JP (1998) Neurofibrillary tangles and Alzheimer's 1029
 966 disease. *Eur Neurol* **40**, 130-140. 1030
- 967 [12] Hasegawa M (2016) Molecular mechanisms in the patho- 1031
 968 genesis of Alzheimer's disease and tauopathies-prion-like 1032
 969 seeded aggregation and phosphorylation. *Biomolecules* **6**, 1033
 970 24. 1034
- 971 [13] Shin WS, Di J, Cao Q, Li B, Seidler PM, Murray KA, Bitan 1035
 972 G, Jiang L (2019). Amyloid β -protein oligomers promote 1036
 973 the uptake of tau fibril seeds potentiating intracellular tau 1037
 974 aggregation. *Alzheimers Res Ther* **11**, 86. 1038
- 975 [14] Wyss-Coray T, Lin C, Yan F, Yu GQ, Rohde M, McConlogue 1039
 976 L, Masliah E, Mucke L (2001) TGF- β 1 promotes 1040
 977 microglial amyloid-beta clearance and reduces plaque bur- 1041
 978 den in transgenic mice. *Nat Med* **7**, 612-618. 1042
- 979 [15] Kanekiyo T, Liu CC, Shinohara M, Li J, Bu G (2012) 1043
 980 LRP1 in brain vascular smooth muscle cells mediates local 1044
 981 clearance of Alzheimer's amyloid- β . *J Neurosci* **32**, 16458- 1045
 982 16465. 1046
- 983 [16] Kanekiyo T, Cirrito JR, Liu CC, Shinohara M, Li J, Schuler 1047
 984 DR, Shinohara M, Holtzman DM, Bu G (2013) Neuronal 1048
 985 clearance of amyloid- β by endocytic receptor LRP1. *J Neuro- 1049
 986 sci* **33**, 19276-19283. 1050
- 987 [17] Zhao Z, Sagare AP, Ma Q, Halliday MR, Kong P, Kisler 1051
 988 K, Winkler EA, Ramanathan A, Kanekiyo T, Bu G, Owens 1052
 989 NC, Rege SV, Si G, Ahuja A, Zhu D, Miller CA, Schneider 1053
 990 JA, Maeda M, Maeda T, Sugawara T, Ichida JK, Zlokovic 1054
 991 BV (2015) Central role for PICALM in amyloid- β blood- 1055
 992 brain barrier transcytosis and clearance. *Nat Neurosci* **18**, 1056
 993 978-987. 1057
- 994 [18] Bell RD, Sagare AP, Friedman AE, Bedi GS, Holtzman 1058
 995 DM, Deane R, Zlokovic BV (2007) Transport pathways for 1059
 996 clearance of human Alzheimer's amyloid beta-peptide and 1060
 997 apolipoproteins E and J in the mouse central nervous system. 1061
 998 *J Cereb Blood Flow Metab* **27**, 909-918. 1062
- 999 [19] Leal MC, Magnani N, Villordo S, Buslje CM, Evelson P, 1063
 1000 Castaño EM, Morelli L (2013) Transcriptional regulation of 1064
 1001 insulin-degrading enzyme modulates mitochondrial amy- 1065
 1002loid β (A β) peptide catabolism and functionality. *J Biol 1066
 1003 Chem* **288**, 12920-12931. 1067
- 1004 [20] Hama E, Shirotani K, Iwata N, Saido TC (2004) Effects of 1068
 1005 neprilysin chimeric proteins targeted to subcellular com- 1069
 1006 partments on amyloid beta peptide clearance in primary 1070
 1007 neurons. *J Biol Chem* **279**, 30259-30264. 1071
- 1008 [21] Lee MJ, Lee JH, Rubinsztein DC (2013) Tau degradation: 1072
 1009 The ubiquitin-proteasome system versus the autophagy- 1073
 1010 lysosome system. *Prog Neurobiol* **105**, 49-59. 1074
- 1011 [22] Subhramanyam CS, Wang C, Hu Q, Dheen ST (2019) 1075
 1012 Microglia-mediated neuroinflammation in neurodegenera- 1076
 1013 tive diseases. *Semin Cell Dev Biol* **94**, 112-120. 1077
- [23] Pereira CF, Santos AE, Moreira PI, Pereira AC, Sousa FJ, 1014
 Cardoso SM, Cruz MT (2019) Is Alzheimer's disease an 1015
 inflammasomopathy? *Ageing Res Rev* **56**, 100966. 1016
- [24] Makin S (2018) The amyloid hypothesis on trial. *Nature* 1017
559, S4-S7. 1018
- [25] Di Paolo G, Kim TW (2011) Linking lipids to Alzheimer's 1019
 disease: Cholesterol and beyond. *Nat Rev Neurosci* **12**, 284- 1020
 296. 1021
- [26] The Gene Ontology Consortium (2019) The Gene Ontology 1022
 Resource: 20 years and still GOing strong. *Nucleic Acids 1023
 Res* **47**, D330-D338. 1024
- [27] Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, 1025
 Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT, 1026
 Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, 1027
 Matese JC, Richardson JE, Ringwald M, Rubin GM, Sher- 1028
 lock G (2000) Gene ontology: Tool for the unification of 1029
 biology. The Gene Ontology Consortium. *Nat Genet* **25**, 1030
 25-29. 1031
- [28] International Society for Biocuration (2018) Biocuration: 1032
 Distilling data into knowledge. *PLoS Biol* **16**, e2002846. 1033
- [29] UniProt Consortium (2019) UniProt: A worldwide hub of 1034
 protein knowledge. *Nucleic Acids Res* **47**, D506-D515. 1035
- [30] Van Cauwenberghe C, Van Broeckhoven C, Sleegers K 1036
 (2016) The genetic landscape of Alzheimer disease: Clinical 1037
 implications and perspectives. *Genet Med* **18**, 421-430. 1038
- [31] Kramarz B, Roncaglia P, Meldal BHM, Huntley RP, Martin 1039
 MJ, Orchard S, Parkinson H, Brough D, Bandopadhyay R, 1040
 Hooper NM, Lovering RC (2018) Improving the gene ontol- 1041
 ogy resource to facilitate more informative analysis and 1042
 interpretation of Alzheimer's disease data. *Genes (Basel)* 1043
9, 593. 1044
- [32] Kramarz B, Huntley RP, Rodríguez-López M, Roncaglia 1045
 P, Saverimuttu SCC, Parkinson H, Bandopadhyay R, Mar- 1046
 tin MJ, Orchard S, Hooper NM, Brough D, Lovering RC 1047
 (2020) Gene ontology curation of neuroinflammation biol- 1048
 ogy improves the interpretation of Alzheimer's disease gene 1049
 expression data. *J Alzheimers Dis*, doi:10.3233/jad-20020 1050
- [33] Mendez D, Gaulton A, Bento AP, Chambers J, De Veij M, 1051
 Félix E, Magariños MP, Mosquera JF, Mutowo P, Nowotka 1052
 M, Gordillo-Marañón M, Hunter F, Junco L, Mugumbate 1053
 G, Rodríguez-Lopez M, Atkinson F, Bosc N, Radoux CJ, 1054
 Segura-Cabrera A, Hersey A, Leach AR (2019) ChEMBL: 1055
 Towards direct deposition of bioassay data. *Nucleic Acids 1056
 Res* **47**, D930-D940. 1057
- [34] Carvalho-Silva D, Pierleoni A, Pignatelli M, Ong C, Fumis 1058
 L, Karamanis N, Carmona M, Faulconbridge A, Hercules 1059
 A, McAuley E, Miranda A, Peat G, Spitzer M, Barrett J, 1060
 Hulcoop DG, Papa E, Koscielny G, Dunham I (2019) Open 1061
 targets platform: New developments and updates two years 1062
 on. *Nucleic Acids Res* **47**, D1056-D1065. 1063
- [35] Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, 1064
 Broackes-Carter F, Campbell NH, Chavali G, Chen C, del- 1065
 Toro N, Duesbury M, Dumousseau M, Galeota E, Hinz 1066
 U, Iannuccelli M, Jagannathan S, Jimenez R, Khadake J, 1067
 Lagreid A, Licata L, Lovering RC, Meldal B, Melidoni AN, 1068
 Milagros M, Peluso D, Perfetto L, Porras P, Raghunath A, 1069
 Ricard-Blum S, Roehert B, Stutz A, Tognolli M, van Roey 1070
 K, Cesareni G, Hermjakob H (2014) The MIntAct project- 1071
 IntAct as a common curation platform for 11 molecular 1072
 interaction databases. *Nucleic Acids Res* **42**, D358-D363. 1073
- [36] Breuza L, Poux S, Estreicher A, Famiglietti ML, Magrane 1074
 M, Tognolli M, Bridge A, Baratin D, Redaschi N; UniProt 1075
 consortium (2016) The UniProtKB guide to the human pro- 1076
 teome. *Database (Oxford)* **2016**, bav120. 1077

- 1078 [37] Meldal BHM, Bye-A-Jee H, Gajdoš L, Hammerová Z, Horácková A, Melicher F, Peretto L, Pokorný D, Lopez MR, Türková A, Wong ED, Xie Z, Casanova EB, Del-Toro N, Koch M, Porras P, Hermjakob H, Orchard S (2019) Complex Portal 2018: Extended content and enhanced visualization tools for macromolecular complexes. *Nucleic Acids Res* **47**, D550-D558.
- 1082 [38] Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T (2003) Cytoscape: A software environment for integrated models of biomolecular interaction networks. *Genome Res* **13**, 2498-2504.
- 1088 [39] Pillich RT, Chen J, Rynkov V, Welker D, Pratt D (2017) NDEx: A community resource for sharing and publishing of biological networks. *Methods Mol Biol* **1558**, 271-301.
- 1092 [40] Bindea G, Mlecnik B, Hackl H, Charoentong P, Tosolini M, Kirilovsky A, Fridman WH, Pagès F, Trajanoski Z, Galon J (2009) ClueGO: A Cytoscape plug-in to decipher functionally grouped gene ontology and pathway annotation networks. *Bioinformatics* **25**, 1091-1093.
- 1097 [41] Orchard S, Kerrien S, Abbani S, Aranda B, Bhate J, Bidwell S, Bridge A, Briganti L, Brinkman FS, Cesareni G, Chatr-aryamontri A, Chautard E, Chen C, Dumousseau M, Goll J, Hancock RE, Hannik LI, Jurisica I, Khadake J, Lynn DJ, Mahadevan U, Peretto L, Raghunath A, Ricard-Blum S, Roechert B, Salwinski L, Stümpflen V, Tyers M, Uetz P, Xenarios I, Hermjakob H (2012) Protein interaction data curation: The International Molecular Exchange (IMEx) consortium. *Nat Methods* **9**, 345-350.
- 1106 [42] Murphy MP, LeVine H 3rd (2010) Alzheimer's disease and the amyloid-beta peptide. *J Alzheimers Dis* **19**, 311-323.
- 1108 [43] Morgat A, Lombardot T, Axelsen KB, Aimo L, Niknejad A, Hyka-Nouspikel N, Coudert E, Pozzato M, Pagni M, Moretti S, Rosanoff S, Onwubiko J, Bougueleret L, Xenarios I, Redaschi N, Bridge A (2017) Updates in Rhea - an expert curated resource of biochemical reactions. *Nucleic Acids Res* **45**, D415-D418.
- 1114 [44] Hastings J, Owen G, Dekker A, Ennis M, Kale N, Muthukrishnan V, Turner S, Swainston N, Mendes P, Steinbeck C (2016) ChEBI in 2016: Improved services and an expanding collection of metabolites. *Nucleic Acids Res* **44**, D1214-D1219.
- 1119 [45] Mielke MM, Lyketsov CG (2010) Alterations of the sphingolipid pathway in Alzheimer's disease: New biomarkers and treatment targets? *Neuromolecular Med* **12**, 331-340.
- 1122 [46] Iwata M, Watanabe S, Yamane A, Miyasaka T, Misonou H (2019) Regulatory mechanisms for the axonal localization of tau protein in neurons. *Mol Biol Cell* **30**, 2441-2457.
- 1125 [47] Arima K (2006) Ultrastructural characteristics of tau filaments in tauopathies: Immuno-electron microscopic demonstration of tau filaments in tauopathies. *Neuropathology* **26**, 475-483.
- 1129 [48] Barbier P, Zejneli O, Martinho M, Lasorsa A, Belle V, Smet-Nocca C, Tsvetkov PO, Devred F, Landrieu I (2019) Role of tau as a microtubule-associated protein: Structural and functional aspects. *Front Aging Neurosci* **11**, 204.
- 1133 [49] Kimura T, Sharma G, Ishiguro K, Hisanaga SI (2018) Phospho-tau bar code: Analysis of phosphoisotypes of tau and its application to tauopathy. *Front Neurosci* **12**, 44.
- 1136 [50] Kimura T, Ishiguro K, Hisanaga SI (2014) Physiological and pathological phosphorylation of tau by Cdk5. *Front Mol Neurosci* **7**, 65.
- 1139 [51] Lee MS, Kwon YT, Li M, Peng J, Friedlander RM, Tsai LH (2000) Neurotoxicity induces cleavage of p35 to p25 by calpain. *Nature* **405**, 360-364.
- 1142 [52] Patrick GN, Zukerberg L, Nikolic M, de la Monte S, Dikkes P, Tsai LH (1999) Conversion of p35 to p25 deregulates Cdk5 activity and promotes neurodegeneration. *Nature* **40**, 615-622.
- 1145 [53] Tanzi RE, Bertram L (2005) Twenty years of the Alzheimer's disease amyloid hypothesis: A genetic perspective. *Cell* **120**, 545-555.
- 1148 [54] Goldman JS, Hahn SE, Catania JW, LaRusse-Eckert S, Rumbaugh M, Strecker MN, Roberts JS, Burke W, Mayeux R, Bird T (2011) Genetic counseling and testing for Alzheimer disease: Joint practice guidelines of the American College of Medical Genetics and the National Society of Genetic Counselors. *Genet Med* **13**, 597-605.
- 1154 [55] Famiglietti ML, Estreicher A, Breuza L, Poux S, Redaschi N, Xenarios I, Bridge A; UniProt Consortium (2019) An enhanced workflow for variant interpretation in UniProtKB/Swiss-Prot improves consistency and reuse in ClinVar. *Database (Oxford)* **2019**, baz040.
- 1159 [56] McGarvey PB, Nightingale A, Luo J, Huang H, Martin MJ, Wu C; UniProt Consortium (2019) UniProt genomic mapping for deciphering functional effects of missense variants. *Hum Mutat* **40**, 694-705.
- 1163 [57] IMEx Consortium Curators, Del-Toro N, Duesbury M, Koch M, Peretto L, Shrivastava A, Ochoa D, Wagih O, Piñero J, Kotlyar M, Pastrello C, Beltrao P, Furlong LI, Jurisica I, Hermjakob H, Hermjakob H, Orchard S, Porras P (2019) Capturing variation impact on molecular interactions in the IMEx Consortium mutations data set. *Nat Commun* **10**, 10.
- 1169 [58] Dana JM, Gutmanas A, Tyagi N, Qi G, O'Donovan C, Martin M, Velankar S (2019) SIFTS: Updated Structure Integration with Function, Taxonomy and Sequences resource allows 40-fold increase in coverage of structure-based annotations for proteins. *Nucleic Acids Res* **47**, D482-D489.
- 1175 [59] Watkins X, Garcia LJ, Pundir S, Martin MJ, UniProt Consortium (2017) ProtVista: Visualization of protein sequence annotations. *Bioinformatics* **33**, 2040-2041.
- 1178 [60] Jeong W, Lee H, Cho S, Seo J (2019) ApoE4-induced cholesterol dysregulation and its brain cell type-specific implications in the pathogenesis of Alzheimer's disease. *Mol Cells* **42**, 739-746.
- 1182 [61] Karch CM, Goate AM (2015) Alzheimer's disease risk genes and mechanisms of disease pathogenesis. *Biol Psychiatry* **77**, 43-51.
- 1185 [62] Andrew RJ, Fisher K, Heesom KJ, Kellett KAB, Hooper NM (2019) Quantitative interaction proteomics reveals differences in the interactomes of amyloid precursor protein isoforms. *J Neurochem* **149**, 399-412.
- 1189 [63] Wang H, Muiznieks LD, Ghosh P, Williams D, Solaris M, Fang A, Ruiz-Riquelme A, Pomès R, Watts JC, Chakrabarty A, Wille H, Sharpe S, Schmitt-Ulms G (2017) Somatostatin binds to the human amyloid β peptide and favors the formation of distinct oligomers. *Elife* **6**, e28401.
- 1194 [64] Clark LN, Poorkaj P, Wszolek Z, Geschwind DH, Nasreddine ZS, Miller B, Li D, Payami H, Awert F, Markopoulou K, Andreadis A, D'Souza I, Lee VM, Reed L, Trojanowski JQ, Zhukareva V, Bird T, Schellenberg G, Wilhelmsen KC (1998) Pathogenic implications of mutations in the tau gene in pallido-ponto-nigral degeneration and related neurodegenerative disorders linked to chromosome 17. *Proc Natl Acad Sci U S A* **95**, 13103-13107.
- 1202 [65] Unawardana CG, Mehrabian M, Wang X, Mueller I, Lubambo IB, Jonkman JEN, Wang H, Schmitt-Ulms G (2015) The human tau interactome: Binding to the

- 1205 ribonucleoproteome, and impaired binding of the proline-
1206 to-leucine mutant at Position 301 (P301L) to chaperones
1207 and the proteasome. *Mol Cell Proteomics* **14**, 3000-3014.
- [66] Perreau VM, Orchard S, Adlard PA, Bellingham SA, Cappa-
1208 pai R, Ciccotosto GD, Cowie TF, Crouch PJ, Duce JA, Evin
1209 G, Faux NG, Hill AF, Hung YH, James SA, Li QX, Mok
1210 SS, Tew DJ, White AR, Bush AI, Hermjakob H, Masters
1211 CL (2010) A domain level interaction network of amy-
1212 loid precursor protein and Abeta of Alzheimer's disease.
1213 *Proteomics* **10**, 2377-2395.
- [67] Villaveces JM, Jiménez RC, Porras P, Del-Toro N, Dues-
1214 bury M, Dumousseau M, Orchard S, Choi H, Ping P, Zong
1215 NC, Askenazi M, Habermann BH, Hermjakob H (2015)
1216 Merging and scoring molecular interactions utilising exist-
1217 ing community standards: Tools, use-cases and a case study.
1218 *Database (Oxford)* **2015**, bau131.
- [68] Hussaini SMQ, Jang MH (2018) New roles for old glue:
1219 Astrocyte function in synaptic plasticity and neurological
1220 disorders. *Int Neurol J* **22**, S106-S114.
- [69] Nagarsheth MH, Viehman A, Lippa SM, Lippa CF (2006)
1221 Notch-1 immunoexpression is increased in Alzheimer's and
1222 Pick's disease. *J Neurol Sci* **244**, 111-116.
- [70] Malhotra A, Younesi E, Sahadevan S, Zimmermann J,
1223 Hofmann-Apitius M (2015) Exploring novel mechanistic
1224 insights in Alzheimer's disease by assessing reliability of
1225 protein interactions. *Sci Rep* **5**, 13634.
- [71] Jupe S, Fabregat A, Hermjakob H (2015) Expression data
1226 analysis with Reactome. *Curr Protoc Bioinformatics* **49**,
1227 8.20.1-8.20.9.
- [72] Orchard S, Hermjakob H (2015) Shared resources, shared
1228 costs—leveraging biocuration resources. *Database (Oxford)*
1229 **2015**, bav009.
- [73] Porras P, Duesbury M, Fabregat A, Ueffing M, Orchard
1230 S, Gloeckner CJ, Hermjakob H (2015) A visual review
1231 of the interactome of LRRK2: Using deep-curated molec-
1232 ular interaction data to represent biology. *Proteomics* **15**,
1233 1390-1404.
- [74] Foulger RE, Denny P, Hardy J, Martin MJ, Sawford T,
1234 Lovering RC (2016) Using the gene ontology to annotate
1235 key players in Parkinson's disease. *Neuroinformatics* **14**,
1236 297-304.
- [75] Pundir S, Magrane M, Martin MJ, O'Donovan C; UniProt
1237 Consortium (2015) Searching and Navigating UniProt
1238 Databases. *Curr Protoc Bioinformatics* **50**, 1.27.1-10.
- [76] Kramarz B, Lovering RC (2019) Gene ontology: A resource
1239 for analysis and interpretation of Alzheimer's disease data.
1240 In *Alzheimer's Disease*. Codon Publications, Brisbane.
- [77] Mlecnik B, Galon J, Bindea G (2019) Automated explora-
1241 tion of gene ontology term and pathway networks with
1242 ClueGO-REST. *Bioinformatics* **35**, 3864-3866.
- [78] Meldal BHM, Orchard S (2018) Searching and extracting
1243 data from the EMBL-EBI complex portal. *Methods Mol Biol*
1244 **1764**, 377-390.