



HAL
open science

Shedding Light on the Ghost Proteome

Tristan Cardon, Isabelle Fournier, Michel Salzet

► **To cite this version:**

Tristan Cardon, Isabelle Fournier, Michel Salzet. Shedding Light on the Ghost Proteome. Trends in Biochemical Sciences, 2021, 46 (3), pp.239-250. 10.1016/j.tibs.2020.10.003 . hal-04446903

HAL Id: hal-04446903

<https://hal.science/hal-04446903v1>

Submitted on 22 Jul 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Shedding light on the Ghost proteome

Tristan Cardon^{1*}, Isabelle Fournier^{1,2*} and Michel Salzet^{1,2*}

¹Univ. Lille, Inserm, CHU Lille, U1192 - Protéomique Réponse Inflammatoire Spectrométrie de Masse - PRISM, F-59000 Lille, France

²Institut Universitaire de France, Paris, France

Abstract

Conventionally, eukaryotic messenger ribonucleic acid (mRNAs) were thought to be monocistronic, leading to the translation of a single protein. However, large-scale proteomics has led to a massive identification of proteins translated from alternative open reading frame (AltORFs) in mRNAs. AltORFs are found in addition to predicted reference ORF (RefORF) or from non-coding RNA (ncRNAs). AltProts are not represented in the conventional protein databases and this “Ghost proteome” was not considered until recently. Some of these proteins are functional and there is growing evidence that they are involved in central functions in physiological and physiopathological context. Here, we review how this ghost proteome fills the gap in our understanding of signaling pathways, establishes new markers of pathologies and highlights therapeutic targets.

Keywords: Alternative Protein (AltProt), smProt, SEPs, small ORF (smORF), mass spectrometry, hidden proteome

* Corresponding authors: Prof. Michel Salzet (ORCID: 0000-0003-4318-0817), Prof. Isabelle Fournier (ORCID: 0000-0003-1096-5044), PhD Tristan Cardon (ORCID: 0000-0003-1751-0528) Laboratoire Protéomique, Réponse Inflammatoire et Spectrométrie de Masse (PRISM) - Inserm U1192 - Université de Lille, Bât SN3, 1er étage, Cité Scientifique, F-59655 Villeneuve d’Ascq Cedex, France. Phone: +33 (0)3 20 43 41 94; Fax: +33 (0)3 20 43 40 54; email: michel.salzet@univ-lille.fr, isabelle.fournier@univ-lille.fr, tristan.cardon@univ-lille.fr

49 **The hidden side of the eukaryotes protein landscape**

50 Proteins are key players of cell regulation and modern Mass Spectrometry (MS) offers
51 the access to their large-scale identification; though large-scale proteomic strategies rely on
52 the comparison of experimental MS data to theoretical prediction through protein database
53 interrogation. Therefore, since the mid 80's, there has been a tremendous effort from the
54 scientific community to provide reliable protein databases [1–3]. Most proteins included in
55 databases are predicted from gene annotations according to well established rules. By
56 convention, for eukaryotes, unlike prokaryotes, annotations only consider genes or transcripts
57 with ORFs longer than 100 codons, shorter products are disregarded as being noncoding
58 RNAs (ncRNAs). Moreover, only a single ORF per coding gene or coding transcript, referred
59 as the reference ORF (RefORF), is considered, since eukaryotes are known to be
60 **monocistronic (see Glossary)** [4,5]. However, it was demonstrated that a significant number
61 of proteins take their origins from alternative ORF (AltORFs) and were missing from
62 databases, making them invisible to the scientific community for a long time [6,7]. Most
63 importantly, these unknown products from AltORF, have been shown to be functionally
64 active. Because they are from different mechanisms and originate from differently classified
65 RNAs, we will use the term Alternative Proteins (AltProts) regardless of their origin and
66 amino acids length (**see Box 1**) [8]. This missing part of the proteome, called the Ghost
67 proteome, is the hidden part of the iceberg. This review will be devoted to these Ghost
68 proteins; reviewing their origin, the approaches used for their identification and to unveil their
69 functions, and the impact they have on the biochemical landscape of the cell, including the
70 pathophysiological functions.

71

72 **Origin of the Ghost proteome**

73 The monocistronic mRNA dogma

74 In eukaryotes, mRNAs carry the genetic information that is translated by ribosomes into
75 protein. Briefly, the main steps of translation are i) the formation of the 43S preinitiation
76 complex, ii) preparation of the mRNA and iii) RNA scanning and the association of the
77 ribosome subunits (**Figure 1A**). These steps are subjected to regulations, like inhibition by
78 protein binding RNA or activation by stimulation of the initiation, resulting in the
79 modification of the protein expression. Regulation can also be affected by modifications of
80 the translation actors (*e.g.* hydrolysis of eIF2, phosphorylation of eIF4E or eIF1, 3 or 5), of
81 the mRNA due to the binding of other proteins at the 5' UTR region causing inaccessibility of
82 the **43S ribosome**. But also, to the micro-RNA binding and blocking the accessibility of
83 mRNA regions to the translation factors [9]. Despite knowledge of these regulation
84 mechanisms, eukaryotic mRNA has always been considered as monocistronic, leading to the
85 translation of a single protein named RefProt. An important amount of work aiming to predict
86 protein sequences was performed by Dr. M. Kozak to understand which sequence is translated
87 into a protein. This work has led Kozak to introduce a scanning model for translation in 1992
88 [4,5,10], where the ribosome is moving along the mRNA sequence from the 5' down to the 3'
89 UTR. In this model, the first "AUG" start codon read is generally the one to initiate
90 translation. Kozak was describing a strong context favoring the initiation of the translation
91 [11], such as the sequence "gccPccAUGG" in which the capital letters are highly conserved, P
92 being a purine (most frequently an adenosine). The strong and weak context described by
93 Kozak explain the possible translation from other codons than the first "AUG" through
94 alternative mechanisms of translation notably by **reinitiation or leaky scanning** [4]. This
95 model has served as structure for the creation of the first protein database, with the
96 UniProtKB/Swiss-Prot database following the same rules for translation. According to these
97 rules, the coding DNA sequence (CDS) region is the longest nucleotide sequence flanked on

98 each side by a start and a stop codon. CDS defines the RefORF which is translated into the
99 RefProt.

100 Existence of mechanisms leading to more than one protein from a single gene is an
101 evolutionary sign that is already established in Bacteria, Archaea, and Viruses. In viruses,
102 multiplicity of AUG initiation sites for translation in the RNA lead to express several proteins
103 in a well-coordinated manner [12]. In eukaryotes, genes which possess multiple initiation site
104 have already been described [13]. They are mostly proteins involved in regulation, such as
105 translation factors, proto-oncogenes, hormone receptors, growth factors, inflammatory
106 mediators, and proteins involved in signal transmission [14–16]. Despite these mechanisms
107 are now considered by prediction algorithms for building databases (e.g. **TrEMBL**) [17] and
108 identify new proteins including their mechanisms; a large part of the proteome has still
109 remained unreferenced.

110

111 Non-coding RNAs translates into Ghost proteins

112 If translation is highly complex, the transcriptional landscape of organisms is as complex.
113 Majority of DNA genomic sequences are transcribed into RNA sequences which are
114 classified between protein-coding RNAs and ncRNAs [18]. Among ncRNAs, the RNA
115 transcripts >200 nucleic acids are called long non-coding RNAs (lncRNAs). They were
116 believed to be molecular relics from the last universal common ancestor of the RNA world
117 and considered as nonfunctional nucleic acid sequences. However, ncRNAs are conserved
118 across most of the cellular life and it was shown that ncRNAs have a vital importance in cells.
119 Although not allowing the translation of proteins, they were demonstrated to be almost
120 uniquely involved in the regulation of the information flow from DNA to protein [19,20].
121 ncRNAs have been shown to play on epigenetic regulation by modulating the accessibility of
122 chromatin by polymerase, on histone methylation [21], and on transcriptome regulation to
123 control protein expression. But they have also been demonstrated to exert other functions,
124 such as assembling an lncRNA with the RNA-binding protein, TLS [22], as a cofactor
125 modulating transcription factors. Involved in, the protein expression [23] or influencing the
126 choice of RNA polymerase II binding site [24]. Finally, lncRNAs also enable post-
127 transcriptional regulation, as described in ZEB2. The fixation of an lncRNA at the 5' UTR
128 region of RNA causes a modification in the RNA splicing with preservation of the fixation
129 site to the ribosome leading to the production of ZEB2 [25]. Interestingly, ZEB2 expression is
130 subjected to several regulations by lncRNAs, including *i.e.* SPRY4-IT1 [26], LINC00461
131 [27], HOTAIRM1 [28]. Function of these lncRNAs is also associated with pathophysiological
132 mechanisms. For example, lncRNA imatinib-upregulated (lncRNA-IUR) is described as a
133 tumor suppressor which inhibits STAT5-CD71 signaling pathway [29].

134 However, it has recently been shown that some of these RNAs considered in first place as
135 non-coding, are translated into proteins [30,31]. Ribosome profiling experiments have shown
136 that ribosomes are able to bind and translate ncRNAs [32–34]. These additional functions of
137 ncRNAs have pushed forwards the development of new databases, including additional
138 protein sequences translated from ncRNA and lncRNA, and has contributed to the emergence
139 of a new field of research, the proteogenomics. Moreover, peptides and proteins encoded by
140 ncRNAs (HOXB-AS3, encoded by lncRNA; FBXW7-185aa, PINT-87aa, and SHPRH-146aa,
141 encoded by circular RNA; and miPEP-200a and miPEP-200b, encoded by primary miRNAs)
142 have been shown to be critical players in cancer development and progression, affecting the
143 altered regulation of glucose metabolism, the epithelial-to-mesenchymal transition (EMT),
144 and the ubiquitination pathway [35]. Thus, clearly establishing the role and function of the
145 proteins translated from ncRNAs in the physiological and physiopathological context can be
146 considered as the novel challenge in clinical and research areas, a major effort is necessary to
147 achieve this goal.

148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167

AltORFs in mRNA translate into Ghost proteins

The second origin of AltProts is from the translation of AltORFs from mature mRNA which translate a RefProt. Alternative translation mechanisms were already foreseen from the work of Kozak, but these were considered then as minor events by comparison to the translation of RefProt. Nevertheless, translation products from the non-coding 5' and 3' UTR, overlapping between CDS and 5' or 3' UTR or reading frame shifts of one or two nucleotides according to the reference CDS (referred as CDS +2 or +3, CDS +1 being the reference CDS) are observed (**Figure 1B, 1C**). For example, AltMRV11 protein is located at the 3' UTR of the mRNA translated from the MRV11 gene. AltMRV11^{EGFP} was cloned and transfected into HeLa cells. AltMRV11^{EGFP} was found to be localized in the nucleus [36] when MRV11 is cytoplasmic, reflecting that AltMRV11 is not a proteoform of MRV11, but represents a completely different sequence to MRV11. Furthermore, these proteins may represent sequence homologies with other RefProts issued from another mRNA or even ncRNA (**Figure 1D**). AltMRV11, presents a sequence homology with BRCA1 interacting protein 3 (BRCA1-IP3) [37] and was demonstrated to interact with BRCA1. AltMRV11 is, therefore, possibly a novel BRCA1 interacting partner that was already identified by yeast two hybrid, but mistakenly rejected as a false positive hit [38].

Ghost proteins: from identification to function

Approaches to AltProt Detection

MS-based proteomic strategies have so far demonstrated to be the most efficient approach to reach massive identification of the Ghost proteome if using a suitable databases, containing predicted sequence of AltProts. Most conventionally, non-targeted large-scale proteomic approaches are performed by **bottom-up**. The usually used **bottom-up shotgun** approaches which involve trypsin digestion of all proteins can be limited for smallest AltProts which do not always present enzymatic cleavage sites and other limitations like the ion charge state of peptides avoiding the identification [39]. This often limits the identification of AltProts because proteins identified on a single peptide match are normally not considered to be of confidence in standard proteomic processing workflow. Using large databases including AltProts have also shown to push the currently used algorithms for protein identification to their limit, in term of FDR calculation and potential false positive identification and redundancy in some peptide sequence based on the AltProt/RefProt homology. Those have been established for a few tens thousands of proteins included in databases but not for hundred thousand ones. It is foreseen that considering AltProts in future studies can push forwards the development of new bio-informatic solutions. Another central aspect of AltProt identification is the sample preparation. Optimized sample preparation protocols, adapted from peptides extraction, using MeOH or boiling water have shown to increase by 30% the number of identified AltProts [40–42]. Top-down proteomics, is another efficient strategy for identifying small proteins such as AltProts. From ovarian cancer tissues, spatially-resolved Top-down proteomics led to identify 15 AltProts which were then back-correlated to the MALDI-MS imaging [43].

Ribosome profiling is an alternative approach to proteomics that was used by different groups to unveil AltProts. Nicholas T. Ingolia and coworkers have used ribosome profiling to demonstrate that noncoding regions of mRNA and also 90% of ncRNA have translation sites capable of producing proteins [33]. Combining Ribosome profiling analysis to the prediction of ribosome binding sites on mRNA, enables to define a score of similarity between prediction and observation, named “**Fragment Length organization similarity score**”

196

197 **(FLOSS)** [33]. The measurement of this score gives additional weight to the probability of
198 protein production in noncoding regions as well as in ncRNAs.

199
200

201 Unravelling the functions of Ghost Proteins

202 Identification of AltProts has revealed a new face of proteomics, challenging the one-
203 gene one-protein dogma. However, although many of previous studies (**see Box2**) have shown
204 the identification of translation sites and AltProts, in mammalian cell lines (HEK, HeLa,
205 NCH82) and tissues (ovarian cancer, glioma, rat spinal cord), and even in extracellular
206 vesicles [44], their functions remain largely under-investigated. As a start, predictions based
207 on AltProt sequences can be achieved. Using homologies found between AltProts and
208 RefProts, domains, families, specific sites, or motif can be identified. This way, more than
209 ~15,000 AltProts were identified with a transmembrane domain or in secreted forms [8]
210 supporting the idea that AltProts are involved in the cellular communication. Domains
211 research via **InterProScan** [45] and **Gene Ontology (GO)**, makes possible to highlight
212 involvement of AltProts in signaling pathways, nucleic acid or protein metabolism, transport
213 or cell organization. Targeted studies, based on screening by CRISPR-Cas9, showed that
214 more than 2,350 AltORFs are involved in modification of cell phenotype. Specifically, co-
215 immunoprecipitation (co-IP) interactome of six lncRNAs has revealed, to name but one, that
216 AltProt derived from lncRNA RP11_469A15.2, is involved in cytochrome C oxidase (COX)
217 complex in the mitochondria [46]. Another way to explore AltProt functions is to search for
218 their RefProt interaction partners [46] and then use this protein-protein interaction (PPI) to
219 place them back into the signaling pathways. This strategy was experimented, first, by re-
220 interrogating publicly available data of cross-linking MS (XL-MS) performed from the nuclei
221 of HeLa cells [47,48]. Among the 1,679 cross-link interactions identified, 292 were shown to
222 involve AltProts. 44 of these AltProts were found to interact with 7 RefProts related to
223 ribonucleoproteins, ribosome subunits and zinc finger protein network. Among these AltProts
224 identified the heterotrimer between the RE/poly(U)-binding/degradation factor 1 (AUF1), the
225 Ribosomal protein 10 (RPL10) and AltATAD2 has been more deeply investigated. Docking
226 models suggest the attachment of AUF1 on the external part of RPL10 and the interaction of
227 AltATAD2 on the RPL10 region interacting with 5S ribosomal RNA as a mechanism of
228 regulation of the ribosome [48]. More recently, new interactome studies by XL-MS,
229 conducted on glioma cell lines in the context of cell transition [49], have pointed out the
230 involvement of AltProts with the ribosome confirming previous observations [48], where
231 AltTRNAU1AP, AltEPHA5 and AltMAP2 are found in interaction with TPM4 protein [49].
232 Some of the identified AltProts and their function are listed in **Supplementary Table 1**. This
233 demonstrates that AltProts are largely involved in replication and translation modulation but
234 also in mitochondria, cell signaling and cell motility through cytoskeleton.

235

236 **Ghost proteins and their function in physiological and physiopathological mechanisms**

237 AltProt regulation dynamics

238 As of yet, only very few antibodies have been raised against AltProts for their detection,
239 hence limiting deep functional analysis of specific AltProt candidates. Large-scale analytical
240 pipelines such as above-mentioned XL-MS have shown to be very useful to gather wealth of
241 information at global level on AltProt functions. Considering that AltProt average size is ~57
242 amino acid [36,50], some assumptions can be made based on their size and the dynamics of
243 their production. For example, AltTRP1 from CDS+3 or AltBING-4 from CDS+2 are
244 described to serve as “cryptic T-cell epitopes” which is of particular importance during
245 thymal selection of T lymphocytes, leading to selection of lymphocytes with antiviral or
246 antitumor activities [51]. Interestingly, it was shown, by studying the abundance of AltProts

247 and their associated RefProts, that AltProt regulation can be independent from that of the
248 RefProt issued from same mRNA. For example, the main *MIEF1* translational product is not
249 the canonical 463 amino acid MID51 protein but the small 70 amino acid alternative MID51
250 protein (AltMID51) which is 2-fold more abundant than MID51 [37] though AltMID51
251 appears to have a shorter lifetime. The re-analysis of data [49] assessing control cells vs. cells
252 under phenotypic transition upon Forskolin treatment (a protein kinase A activator inducing
253 cell differentiation and EMT), indicates major abundance variations between AltProts and
254 RefProts from same mRNA. For AltANKZF1 and AltVPS52 (**Figure 2**), AltANKZF1 is
255 found to be 10-fold more abundant than ANKZF1 in the control; however, after Forskolin
256 stimulation, there is an increase of ANKZF1 which reaches then to the same abundance as
257 AltANKZF1. Similar trend is observed with AltVPS52 and VPS52. In the control, AltVPS52
258 is 10-fold more abundant than VPS52 and after Forskolin treatment, AltVPS52 level remains
259 higher than that of VPS52. This reflects the dynamic of AltProt expression, suggesting their
260 potential implication in different signaling pathways. However, considering size similarities
261 between AltProts with a mean length of 50 amino acids and polypeptides or neuropeptides
262 typically 3 to 100 amino acids [52] (like the neuropeptide Y of 36 amino acids [53]) ; it can be
263 speculate that AltProts can show similar functions to these peptides. For example, it is well
264 known that snake, spider and Conus venoms contain venom peptides which are known to
265 interact with receptors such as nicotinic acetylcholine receptors and can interfere with
266 cholinergic transmission or have an impact on cardiac myocyte [54]. Neuropeptides also
267 present a comparable size than some of AltProts. Hence, it is conceivable that among other
268 functions, AltProts could be linked to the activation of orphan receptors.

269

270 AltProts in physiopathological mechanisms

271 AltProts have been described in various pathological contexts among which cancer
272 (**Figure 3**). In ovarian cancer, 15 AltProts have been highlighted, including 4 derived from
273 tumor region in ovarian cancer: AltCMBL, AltGNL1, AltRP11-576E20.1, AltCSNK1A1L
274 [43]. In breast cancer, AltProt CASIMO1, derived from an lncRNA, is described as an
275 oncogene involved in MAPK activation signaling pathway [55]. AltEDARADD which is now
276 annotated by UniProt (accession number L0R849), is shown to be with a high level of
277 expression and associate with a poor prognosis for patients suffering from high grade serous
278 ovarian cancer [40]. In esophageal cancer, lncRNA LINC00278 is described to promote the
279 expression of the YY1BM AltProt [56]. Some AltProts have been as well described in the
280 activation of the EMT [49,57]. AltProt designated as Humanin (HN) polypeptide, was found
281 in Alzheimer disease as a suppressor of genes involved in neuronal death [58]. These
282 evidences demonstrate that in many cases, AltProts are acting in regulation, playing for
283 example on translation of RNAs and therefore regulating gene expression (**Figure 4.A**).
284 AltProts are involved in mechanisms that were thought to be well understood, though
285 apparently one piece of the puzzle was missing. Cancer development is often driven by gene
286 mutations; among others mutations in P53 or BRCA1/BRCA2 RefProts are well known. It is
287 very likely, though not yet demonstrated, that mutations on AltProts are also driving some
288 pathologies (**Figure 4.B**). AltGNL1 which results from a shift in the ORF of the mRNA
289 (nucleic acids 1487-1679) coding for RefProt GNL1 is an excellent example. The COSMIC
290 database [59], which references the mutations of various cancers, describes 15 mutations in
291 the region of the mRNA coding both for GNL1 and AltGNL1. Among the 15 mutations
292 expected for AltGNL1; 12 will modify an amino acid of the protein (**Supplementary Table**
293 **2**). The main one, which is a mutation of the nucleic acid 1646 (C>T, COSM8898738)
294 impacts GNL1 by changing A>V. For AltGNL1, this mutation leads to the appearance of a
295 premature stop codon. Thus, upon the effect of this mutation, AltGNL1 (64 amino acids)
296 translates into a truncated proteoform of 53 amino acids, accordingly to the observations of

297 mRNA expression in endometrium carcinoma. Such modifications can significantly change
298 their response and the cell phenotype. Therefore, tracking for mutations by co-analyses of
299 non-coding regions of mRNA by RNA-seq together with AltProts identification will open the
300 door into new target in cancer induced by mutation.

301

302 **Concluding Remarks**

303 As a result of the rules defining the translation to predict proteins in eukaryotes, the entire
304 Ghost proteome has until recently, remained unknown. MS-based large-scale analysis and
305 identification of ribosome binding sites on RNA by ribosome profiling has revealed its
306 existence. It is now clear that, as in prokaryotes, eukaryote genes are polycistronic. This idea
307 challenges the paradigm in the view of proteomics, transcriptomics and genetics; though, not
308 radically a new idea since Kozak had already described this principle in 1999 [5]. The
309 polycistronic nature of RNA limits the necessary material for synthesizing proteins, hence the
310 production of AltProt is probably an evolutionary sign already carried out by bacteria [60],
311 plants [61] and even some insects [62]. However, this “Lost World” of proteins remains to be
312 thoroughly explored and might help explaining mechanisms that are still poorly understood.
313 For example, one could hypothesize that AltProts are ligands of orphan receptors. But
314 validating such an assumption will require functional analysis and unveiling the function of
315 AltProts is definitely the next challenge. The ability of AltProts to present post-translational
316 modifications [8] and mutations, establish protein-protein interactions [63], be active player in
317 cellular communication. This can be through by extra cellular vesicle release and bind to
318 membranes or cytoskeleton, demonstrates that they share a similar diversity and complexity in
319 their functions than RefProts. Global large-scale analyses is one way to rapidly predict
320 AltProts function by identifying domains and finding their interaction partners. However,
321 studies have only yet lifted the veil on this new world of proteins and there are still many
322 opened questions to be answered (see “**Outstanding Questions**”) and future challenges
323 ahead. Among others the role of AltProts in physiological and pathophysiological
324 mechanisms is one [64]. Another is to grasp knowledge about AltProt mutations in relation to
325 pathologies. Moreover, it is clear that AltProts can potentially interfere with inhibitors of
326 RefProt drug targets which would explain limited efficiency of certain drugs. Further,
327 potentiality of a part of the Ghost proteome to act as a rescue system for cell is a reasonable
328 assumption and something to be investigated. Finally, presence of a subset of mRNA products
329 also rises questions on the strategies developed to target protein expression at DNA level. In
330 particular, CRISPR-Cas9 technology shows a high frequency of off-target activity ($\geq 50\%$)
331 [65,66] which is a major concern to therapeutic and clinical applications [67]. In light of
332 alternative translation mechanisms, such off-targets can lead to the presence of several
333 AltProts with various specific biological functions. The downstream impact of releasing such
334 sub-products is unknown and might induces aberrant mechanisms and an opposite effect. In a
335 nutshell, by unlocking the Ghost proteome, researchers have opened the Pandora’s Box.

336

337 **ACKNOWLEDGMENT**

338 This research was supported by funding from Ministère de l’Enseignement Supérieur, de la
339 Recherche et de l’Innovation (MESRI), Institut National de la Santé et de la Recherche
340 Médicale (Inserm) and Université de Lille.

341 Figures are realized on the Adobe Illustrator software.

342

343 1 Breuza, L. *et al.* (2016) The UniProtKB guide to the human proteome. *Database* 2016,
344 bav120

345 2 Bateman, A. *et al.* (2017) UniProt: The universal protein knowledgebase. *Nucleic Acids*
346 *Res.* 45, D158–D169

347 3 Bateman, A. *et al.* (2015) UniProt: A hub for protein information. *Nucleic Acids Res.*
348 43, D204–D212

349 4 Kozak, M. Regulation of translation in eukaryotic systems. , *Annual Review of Cell*
350 *Biology*, 8. 28-Nov-(1992) , Annual Reviews 4139 El Camino Way, P.O. Box 10139,
351 Palo Alto, CA 94303-0139, USA, 197–225

352 5 Kozak, M. (1999) Initiation of translation in prokaryotes and eukaryotes. *Gene* 234,
353 187–208

354 6 Vanderperre, B. *et al.* (2012) HAltORF: A database of predicted out-of-frame
355 alternative open reading frames in human. *Database* 2012, bas025

356 7 Moulleron, H. *et al.* Death of a dogma: Eukaryotic mRNAs can code for more than
357 one protein. , *Nucleic Acids Research*, 44. 08-Jan-(2016) , Oxford University Press,
358 14–23

359 8 Samandi, S. *et al.* (2017) Deep transcriptome annotation enables the discovery and
360 functional characterization of cryptic small proteins. *Elife* 6, e27860

361 9 Hershey, J.W.B. (1991) Translational Control in Mammalian Cells. *Annu. Rev.*
362 *Biochem.* 60, 717–755

363 10 Kozak, M. (1987) An analysis of 5'-noncoding sequences from 699 vertebrate
364 messenger rNAS. *Nucleic Acids Res.* 15, 8125–8148

365 11 Kozak, M. (1986) Point mutations define a sequence flanking the AUG initiator codon
366 that modulates translation by eukaryotic ribosomes. *Cell* 44, 283–292

367 12 Fouillot, N. *et al.* (1993) Translation of the hepatitis B virus P gene by ribosomal
368 scanning as an alternative to internal initiation. *J. Virol.* 67, 4886–4895

369 13 Kozak, M. An analysis of vertebrate mRNA sequences: Intimations of translational
370 control. , *Journal of Cell Biology*, 115. (1991) , The Rockefeller University Press, 887–
371 903

372 14 Prats, H. *et al.* (1989) High molecular mass forms of basic fibroblast growth factor are
373 initiated by alternative CUG codons. *Proc. Natl. Acad. Sci. U. S. A.* 86, 1836–1840

374 15 Bugler, B. *et al.* (1991) Alternative initiation of translation determines cytoplasmic or
375 nuclear localization of basic fibroblast growth factor. *Mol. Cell. Biol.* 11, 573–577

376 16 Hann, S.R. *et al.* (1988) A non-AUG translational initiation in c-myc exon 1 generates
377 an N-terminally distinct protein whose synthesis is disrupted in Burkitt's lymphomas.
378 *Cell* 52, 185–195

379 17 Möller, S. *et al.* (1999) EDITtoTrEMBL: a distributed approach to high-quality
380 automated protein sequence annotation. *Bioinformatics* 15, 219–27

381 18 Quinn, J.J. and Chang, H.Y. Unique features of long non-coding RNA biogenesis and
382 function. , *Nature Reviews Genetics*, 17. 01-Jan-(2016) , Nature Publishing Group, 47–
383 62

384 19 Poole, A. *et al.* (1999) Early evolution: Prokaryotes, the new kids on the block.
385 *BioEssays* 21, 880–889

386 20 Jeffares, D.C. *et al.* (1998) Relics from the RNA world. *J. Mol. Evol.* 46, 18–36

387 21 Rinn, J.L. *et al.* (2007) Functional Demarcation of Active and Silent Chromatin
388 Domains in Human HOX Loci by Noncoding RNAs. *Cell* 129, 1311–1323

389 22 Wang, X. *et al.* (2008) Induced ncRNAs allosterically modify RNA-binding proteins in
390 cis to inhibit transcription. *Nature* 454, 126–130

391 23 Feng, J. *et al.* (2006) The Evf-2 noncoding RNA is transcribed from the Dlx-5/6
392 ultraconserved region and functions as a Dlx-2 transcriptional coactivator. *Genes Dev.*

- 393 20, 1470–1484
- 394 24 Martianov, I. *et al.* (2007) Repression of the human dihydrofolate reductase gene by a
395 non-coding interfering transcript. *Nature* 445, 666–670
- 396 25 Beltran, M. *et al.* (2008) A natural antisense transcript regulates Zeb2/Sip1 gene
397 expression during Snail1-induced epithelial-mesenchymal transition. *Genes Dev.* 22,
398 756–769
- 399 26 Yao, H. *et al.* (2020) LncRNA SPRY4-IT1 promotes progression of osteosarcoma by
400 regulating ZEB1 and ZEB2 expression through sponging of miR-101 activity. *Int. J.*
401 *Oncol.* 56, 85–100
- 402 27 Li, X. *et al.* (2020) The Lnc LINC00461/miR-30a-5p facilitates progression and
403 malignancy in non-small cell lung cancer via regulating ZEB2. *Cell Cycle* 19, 825–836
- 404 28 Lin, Y.H. *et al.* (2020) Long non-coding RNA HOTAIRM1 promotes proliferation and
405 inhibits apoptosis of glioma cells by regulating the miR-873-5p/ZEB2 axis. *Chin. Med.*
406 *J. (Engl).* 133, 174–182
- 407 29 Wang, X. *et al.* (2019) Novel lncRNA-IUR suppresses Bcr-Abl-induced tumorigenesis
408 through regulation of STAT5-CD71 pathway. *Mol. Cancer* 18, 84
- 409 30 Rion, N. and Rüegg, M.A. LncRNA-encoded peptides: More than translational noise? ,
410 *Cell Research*, 27. 01-May-(2017) , Nature Publishing Group, 604–605
- 411 31 Matsumoto, A. *et al.* (2017) MTORC1 and muscle regeneration are regulated by the
412 LINC00961-encoded SPAR polypeptide. *Nature* 541, 228–232
- 413 32 Ingolia, N.T. *et al.* (2011) Ribosome profiling of mouse embryonic stem cells reveals
414 the complexity and dynamics of mammalian proteomes. *Cell* 147, 789–802
- 415 33 Ingolia, N.T.T. *et al.* (2014) Ribosome Profiling Reveals Pervasive Translation Outside
416 of Annotated Protein-Coding Genes. *Cell Rep.* 8, 1365–1379
- 417 34 Bazzini, A.A. *et al.* (2014) Identification of small ORFs in vertebrates using ribosome
418 footprinting and evolutionary conservation. *EMBO J.* 33, 981–993
- 419 35 Wang, J. *et al.* ncRNA-Encoded Peptides or Proteins and Cancer. , *Molecular Therapy*,
420 27. 02-Oct-(2019) , Cell Press, 1718–1725
- 421 36 Vanderperre, B. *et al.* (2013) Direct Detection of Alternative Open Reading Frames
422 Translation Products in Human Significantly Expands the Proteome. *PLoS One* 8,
423 e70698
- 424 37 Delcourt, V. *et al.* (2018) The Protein Coded by a Short Open Reading Frame, Not by
425 the Annotated Coding Sequence, Is the Main Gene Product of the Dual-Coding Gene
426 MIEF1. *Mol. Cell. Proteomics* 17, 2402–2411
- 427 38 Liu, Y. *et al.* (2011) Yeast two-hybrid junk sequences contain selected linear motifs.
428 *Nucleic Acids Res.* 39, e128
- 429 39 Slavoff, S.A. *et al.* (2013) Peptidomic discovery of short open reading frame-encoded
430 peptides in human cells. *Nat. Chem. Biol.* 9, 59–64
- 431 40 Cardon, T. *et al.* (2019) Optimized Sample Preparation Workflow for Improved
432 Identification of Ghost Proteins. DOI: 10.1021/acs.analchem.9b04188
- 433 41 Cassidy, L. *et al.* (2019) Depletion of High-Molecular-Mass Proteins for the
434 Identification of Small Proteins and Short Open Reading Frame Encoded Peptides in
435 Cellular Proteomes. *J. Proteome Res.* 18, 1725–1734
- 436 42 Cassidy, L. *et al.* (2016) Combination of Bottom-up 2D-LC-MS and Semi-top-down
437 GelFree-LC-MS Enhances Coverage of Proteome and Low Molecular Weight Short
438 Open Reading Frame Encoded Peptides of the Archaeon *Methanosarcina mazei*. *J.*
439 *Proteome Res.* 15, 3773–3783
- 440 43 Delcourt, V. *et al.* (2017) Combined Mass Spectrometry Imaging and Top-down
441 Microproteomics Reveals Evidence of a Hidden Proteome in Ovarian Cancer.
442 *EBioMedicine* 21, 55–64

443 44 Murgoci, A.N. *et al.* (2020) Reference and Ghost Proteins Identification in Rat C6
444 Glioma Extracellular Vesicles. *iScience* 23, 101045

445 45 Jones, P. *et al.* (2014) InterProScan 5: Genome-scale protein function classification.
446 *Bioinformatics* 30, 1236–1240

447 46 Chen, J. *et al.* (2020) Pervasive functional translation of noncanonical human open
448 reading frames. *Science* (80-.). 367, 140–146

449 47 Liu, F. *et al.* (2015) Proteome-wide profiling of protein assemblies by cross-linking
450 mass spectrometry. *Nat. Methods* 12, 1179–1184

451 48 Cardon, T. *et al.* (2019) Nuclei of HeLa cells interactomes unravel a network of ghost
452 proteins involved in proteins translation. *Biochim. Biophys. Acta - Gen. Subj.* DOI:
453 10.1016/j.bbagen.2019.05.009

454 49 Cardon, T. *et al.* (2020) Alternative proteins are functional regulators in cell
455 reprogramming by PKA activation. *Nucleic Acids Res.* DOI: 10.1093/nar/gkaa277

456 50 Delcourt, V. *et al.* (2018) Small Proteins Encoded by Unannotated ORFs are Rising
457 Stars of the Proteome, Confirming Shortcomings in Genome Annotations and Current
458 Vision of an mRNA. *Proteomics* 18, 1700058

459 51 Ho, O. and Green, W.R. (2006) Unexpected Cryptic T Cell Epitopes: Expecting the
460 Alternative Translational Products and. DOI: 10.4049/jimmunol.177.12.8283

461 52 Wang, Y. *et al.* (2015) NeuroPep: A comprehensive resource of neuropeptides.
462 *Database* 2015,

463 53 Reichmann, F. and Holzer, P. Neuropeptide Y: A stressful review. , *Neuropeptides*, 55.
464 01-Feb-(2016) , Churchill Livingstone, 99–109

465 54 Utkin, Y.N. (2019) Last decade update for three-finger toxins: Newly emerging
466 structures and biological activities. *World J. Biol. Chem.* 10, 17–27

467 55 Polycarpou-Schwarz, M. *et al.* (2018) The cancer-associated microprotein CASIMO1
468 controls cell proliferation and interacts with squalene epoxidase modulating lipid
469 droplet formation. *Oncogene* 37, 4750–4768

470 56 Wu, S. *et al.* (2020) A Novel Micropeptide Encoded by Y-Linked LINC00278 Links
471 Cigarette Smoking and AR Signaling in Male Esophageal Squamous Cell Carcinoma.
472 *Cancer Res.* DOI: 10.1158/0008-5472.can-19-3440

473 57 Vergara, D. *et al.* (2020) A Hidden Human Proteome Signature Characterizes the
474 Epithelial Mesenchymal Transition Program. *Curr. Pharm. Des.* 26, 372–375

475 58 Hashimoto, Y. *et al.* (2001) A rescue factor abolishing neuronal cell death by a wide
476 spectrum of familial Alzheimer’s disease genes and A β . *Proc. Natl. Acad. Sci. U. S.*
477 *A.* 98, 6336–41

478 59 Bamford, S. *et al.* (2004) The COSMIC (Catalogue of Somatic Mutations in Cancer)
479 database and website. *Br. J. Cancer* 91, 355–358

480 60 Güell, M. *et al.* Bacterial transcriptomics: What is beyond the RNA horiz-ome? ,
481 *Nature Reviews Microbiology*, 9. Sep-(2011) , 658–669

482 61 Hanada, K. *et al.* Small open reading frames associated with morphogenesis are hidden
483 in plant genomes. DOI: 10.1073/pnas.1213958110

484 62 Ladoukakis, E. *et al.* (2011) Hundreds of putatively functional small open reading
485 frames in *Drosophila*. *Genome Biol.* 12, R118

486 63 Cardon, T. *et al.* (2020) Optimized Sample Preparation Workflow for Improved
487 Identification of Ghost Proteins. *Anal. Chem.* 92, 1122–1129

488 64 Le Rhun, E. *et al.* (2017) Evaluation of non-supervised MALDI mass spectrometry
489 imaging combined with microproteomics for glioma grade III classification. *Biochim.*
490 *Biophys. Acta - Proteins Proteomics* 1865, 875–890

491 65 Fu, Y. *et al.* (2013) High-frequency off-target mutagenesis induced by CRISPR-Cas
492 nucleases in human cells. *Nat. Biotechnol.* 31, 822–826

493 66 Tuladhar, R. *et al.* (2019) CRISPR-Cas9-based mutagenesis frequently provokes on-
494 target mRNA misregulation. *Nat. Commun.* 10, 4056

495 67 Han, H.A. *et al.* Mitigating off-target effects in CRISPR/Cas9-mediated in vivo gene
496 editing. , *Journal of Molecular Medicine*, 98. 01-May-(2020) , Springer, 615–632

497 68 Slavoff, S.A. *et al.* (2014) A human short open reading frame (sORF)-Encoded
498 polypeptide that stimulates DNA end joining. *J. Biol. Chem.* 289, 10950–10957

499 69 Hao, Y. *et al.* (2018) SmProt: a database of small proteins encoded by annotated coding
500 and non-coding RNA loci. *Brief. Bioinform.* 19, 636–643

501 70 Peeters, M.K.R. and Menschaert, G. (2020) The hunt for sORFs: A multidisciplinary
502 strategy. *Exp. Cell Res.* DOI: 10.1016/j.yexcr.2020.111923

503 71 Brunet, M.A. *et al.* (2019) OpenProt: A more comprehensive guide to explore
504 eukaryotic coding potential and proteomes. *Nucleic Acids Res.* 47, D403–D410

505 72 Olexiuk, V. *et al.* (2017) An update on sORFs.org: a repository of small ORFs
506 identified by ribosome profiling. *Nucleic Acids Res.* 46, 497–502

507 73 Brunet, M.A. *et al.* (2020) Reconsidering proteomic diversity with functional
508 investigation of small ORFs and alternative ORFs. *Exp. Cell Res.* 393, 112057

509 74 Fairman, R. *et al.* (1993) Multiple oligomeric states regulate the DNA binding of helix-
510 loop-helix peptides. *Proc. Natl. Acad. Sci. U. S. A.* 90, 10429–10433

511 75 Tabb-Massey, A. *et al.* Ribosomal proteins Rps0 and Rps21 of *Saccharomyces*
512 *cerevisiae* have overlapping functions in the maturation of the 3' end of 18S rRNA.
513 DOI: 10.1093/nar/gkg899

514 76 Vattem, K.M. and Wek, R.C. (2004) Reinitiation involving upstream ORFs regulates
515 ATF4 mRNA translation in mammalian cells. *Proc. Natl. Acad. Sci. U. S. A.* 101,
516 11269–11274

517 77 Chipuk, J.E. and Green, D.R. How do BCL-2 proteins induce mitochondrial outer
518 membrane permeabilization? , *Trends in Cell Biology*, 18. 01-Apr-(2008) , Elsevier
519 Current Trends, 157–164

520 78 O'Leary, N.A. *et al.* (2016) Reference sequence (RefSeq) database at NCBI: Current
521 status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* 44, D733–
522 D745

523 79 Pueyo, J.I. *et al.* New Peptides Under the s(ORF)ace of the Genome. , *Trends in*
524 *Biochemical Sciences*, 41. Aug-(2016) , 665–678

525 80 Ingolia, N.T. Ribosome profiling: New views of translation, from single codons to
526 genome scale. , *Nature Reviews Genetics*, 15. 28-Jan-(2014) , Nature Research, 205–
527 213

528 81 Ingolia, N.T. *et al.* (2009) Genome-wide analysis in vivo of translation with nucleotide
529 resolution using ribosome profiling. *Science (80-)*. 324, 218–223

530

531

532 **Glossary**

533 **43S ribosome:** The preinitiation complex (43S PIC) is a ribonucleoprotein complex that
534 exists during an early step of the translation. The formation of the 43s ribosome complex
535 arises from the recycling of the small ribosome 40S ribosomal subunit bound by the initiation
536 factors (eIF): eIF1, eIF3 and eIF1A. This is followed by the docking of eIF5 and eIF2, itself
537 transporting the tRNA-MET, allowing the recognition of the START codon and the binding
538 of the first amino acid (methionine).

539 **Bottom-up:** Bottom-up is the most used large-scale proteomic approach for the identification
540 and relative quantification of proteins in a single experiment. It involves a digestion of the
541 proteins using one (generally trypsin) or several enzymes (*e.g.* LysC/trypsin) and the MS
542 identification of the proteins through measuring their digestion peptides.

543

544 **Bottom-up shotgun:** Shotgun is a sub-type of bottom-up approaches where the digestion of
545 all proteins together is performed after the extraction; the tryptic peptides being further
546 analyzed by LC-MS. In shotgun the protein identification is solely based on the molecular
547 weight of the peptides and the structural information collected in tandem MS spectra after gas
548 phase fragmentation of the ions.

549

550 **Darkome:** The Darkome designates transcripts, proteins and metabolites that are missing
551 annotation in the databases. For example, the AltProts and their transcripts showing mutations
552 that could be pathology drivers, such as for cancer, as still missing annotation. Indeed,
553 mutations which are considered as silent mutations because not in the region of the transcript
554 sequence related to the RefProt translation could be non-silencing for the AltProts.
555 Altogether, the Darkome represents still hundred thousand of unknown products.

556

557 **Fragment Length Organization Similarity Score (FLOSS):** The FLOSS is defined as,

558
$$0.5 \times \sum_{l=26}^{34} f(l) - f_{ref}(l)$$

559 where $f(l)$ is the fraction of reads of length l in the transcript histogram and $f_{ref}(l)$ is the
560 corresponding fraction in the reference histogram (raw count excluding gene overlap and
561 noncoding).

562

563 **Gene Ontology (GO):** The GO is the largest world knowledgebase source of information,
564 across all species, on the functions of genes that is both human-readable and machine-
565 readable.

566

567 **InterProScan:** InterProScan is a free software that provides functional analysis of proteins by
568 classifying them into families and predicting domains and important sites. It uses predictive
569 models and known signatures from several different databases from the InterPro consortium.

570

571 **Monocistronic:** refers to the ability of an mRNA to be translated into a single protein chain.
572 Eukaryotes are considered monocistronic.

573 **Reinitiation and leaky scanning:** These are 2 mechanisms explaining the translation of
574 AltORFs. The reinitiation mechanism is due to the formation of the 80S ribosome upstream of
575 the most likely AUG site. This small region, formerly called upstream ORF (upORF) is
576 described as not exceeding 30 codons, if its length is greater (more than 120 codons) it no
577 longer allows linking with the main ORF. The upORF region allows the regulation of the
578 transcription of the downstream protein (RefProt), however the role of the small protein is
579 generally attributed to a non-functional peptide. The leaky scanning allows to pass the rule of
580 the first codon "AUG". The second or third Start codon will then be used. The mechanism is
581 due to the absence of an optimal initiation context. Indeed, if the Kozak context around the
582 AUG codon is not strong enough, the 40S subunit does not stop at the expected Start codon
583 and translates the following region.

584

585 **TrEMBL:** Translation of the EMBL (TrEMBL) is a large protein database in SwissProt
586 format generated by computer translation of all coding sequences in the EMBL Nucleotide
587 Sequence Database database, that are not in SWISS-PROT. Unlike SWISS-PROT, TrEMBL
588 entries are manually annotated.

589

590 **Text Boxes:**

591 **Box1: The Alternative Protein kingdom.**

592 The alternative proteins terminology encompasses a large number of proteins issued from
593 alternative mechanisms (e.g. **reinitiation and leaky scanning**) and various origins (mRNA
594 and ncRNA). These translation products give rise to proteins (as experimentally evidenced)
595 which are in average of small size (<57 amino acids). Currently, there is a clear lack of
596 consensus regarding the terminology of this novel class of proteins in the literature, one
597 reason being the different approaches and angles by which they were initially studied by
598 different groups. Indeed, in the literature, they can be referred to as, short ORFs (sORFs)
599 originating *i.e.* sORF-Encoded Polypeptides (SEPs) [68], Small Proteins (smPROT) [69],
600 polypeptides [70] or AltProts [71]. Indeed, sORFs from the sORFs.org repository [72], are
601 described as ORF in the range of 10 to 100 codons, issued from ribosome profiling analysis,
602 therefore accepting other than AUG start codon, but not including shifts in the CDS. While,
603 AltORF coding AltProt are described as >30 codons, starting with an AUG codon but also
604 found in reading frame shifts [73]. In 2012, 17,096 AltProts sequences were predicted from
605 ~31,422 mRNAs or 8,744 genes from the human genome and referenced as the “Human
606 alternative Open Reading Frame” (HaltORF) database [6]. In this initial database, 83% of
607 predicted AltProts were in the CDS+2 and 17% in the CDS+3 (CDS+1 corresponding to the
608 RefProt). However, this was only taking into account the reading offsets, without referencing
609 the predicted translation products in the 5', 3' UTR regions or ncRNA. In 2019, all of the
610 predictions were included and are now available in the public "OpenProt" database [71]. It is
611 currently predicted ~450,000 proteins in human, without considering post-translational
612 modifications. AltProt prediction from other species such as mice, rat and zebrafish are also
613 now accessible in OpenProt. Because we wish to keep here the big picture on these novel
614 proteins independently of their length in amino acids or nucleotides, their origin (mRNA,
615 ncRNA) or their initiation by an AUG start codon, we prefer to refer to them as proteins
616 issued from alternative mechanisms to the RefProts, or as Alternative Proteins (AltProts).

617 **Box2: History of the Ghost proteome evidence**

618 Despite the general dogma on the monocistronic nature of eukaryotes mRNAs, it was already
619 demonstrated 20 years ago that mRNA as well as certain ncRNA could translate into several
620 proteins. Indeed, in 2000 the assessment of the ribosome profiling has led to prove the
621 capacity of the ribosome to bind on noncoding region of mRNAs and ncRNAs [32,33].
622 Before the advent of large-scale proteomics, few proteins issued from other ORFs than the
623 RefORF have been punctually discovered and identified such as id peptide [74], RPS21 [75],
624 AltATF4 [76] or BCL-2 [77]. Since protein databases used in large scale proteomic
625 approaches (e.g. UniprotKB [3] or NCBI [78]) are based on genomic annotations, which rule
626 out AltORFs and ncRNA, AltProts have been missing and hence, could not be identified
627 through MS-based proteomics. More recently, it was demonstrated that proteins issued from
628 AltORFs were far more common than initially expected. The first massive identification of
629 AltProts was obtained in 2011 at the protein level by shotgun proteomics using high-
630 resolution (HR) LC-MS. Surprisingly, large scale proteomic bottom-up experiments were
631 always showing more than 10% of tandem MS data remaining unmatched after protein
632 database interrogation. This has led to consider alternative translation mechanisms in the
633 protein predictions and to generate a first database of human AltProts (HaltORF) [6]; though
634 this was far from being complete. A total of 1,259 AltProts were initially identified by
635 shotgun proteomics from human cell lines, FFPE tissues and fluids, both from normal and
636 pathological samples [43]. These results have opened the door to several *in silico* studies to

637 assess these AltProt databases, leading to the protein interrogation in large-scale proteomic
638 experiments and further evaluate their physiological and pathophysiological roles.

639 In the following years (**Figure I**), high-throughput genome and transcriptome sequencing was
640 used to validate an important number of sORFs and AltORF containing less than 100 codons
641 that were previously arbitrary considered as noncoding sequences. The translation of these
642 sORFs into peptides or small proteins (AltProts, microproteins, SEPs) was demonstrated by
643 different strategies such as ribosome profiling [70,79–81] and/or peptidomics coupled to
644 massively parallel RNA-seq [39]. This accumulating evidence both demonstrated by genome
645 annotation and protein translation, has largely questioned the validity of the gene-protein
646 dogma pushing forward the development of new databases including the AltProts as a
647 mandatory step for their identification in large-scale proteomic studies.

648

649

650 Figure Legends:

651

652 **Figure 1.**

653 **Scheme of AltProts expression regions.** **A)** The conventional vision of the eukaryotic
654 protein translation: where the RefProt expression comes from the first START codon "AUG",
655 with a CDS framed by a 5' & 3'UTR region. **B)** The 5' & 3' UTR regions can also be the
656 location of the ribosome binding, allowing the translation of AltProt^{5'UTR} or AltProt^{3'UTR}. **C)**
657 The ribosome can also shift of the reading frame with a fixation at +2 or +3 of the reference
658 START codon, thus allowing the translation of AltProt^{CDS+2} or AltProt^{CDS+3}. **D)** The last
659 AltProt translation source is ncRNA and lncRNA.

660

661 **Figure 2.**

662 **Variation of the relative abundance (LFQ) of AltProt vs. RefProt** associated with a same
663 mRNA. Each condition has been performed in triplicates **A)** AltANKZF1 vs. ANKZF1. The
664 RefProt ANKZF1 shows a lower relative abundance than its AltProt in the control condition,
665 while they show the same abundance under FRSK stimulation, meaning that FRSK
666 stimulation increase the level of ANKZF1 but does not influence AltANKZF1. **B)** AltVPS52
667 vs. VPS52. AltVPS52 has a higher abundance than the reference VPS52 but both protein
668 show the same regulation trend under FRSK stimulation.

669

670 **Figure 3.**

671 **Pathophysiological implication and identification of AltProt.** The Ghost proteome drives
672 different proteins involved in the development of cancer. AltProts are plotted back on cancer
673 hallmarks showing that already identified AltProt covers most of hallmarks (in blue circle)
674 and cancers (in red square).

675

676 **Figure 4.**

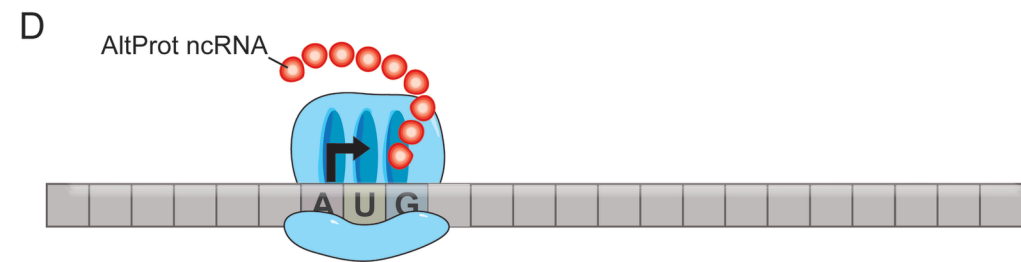
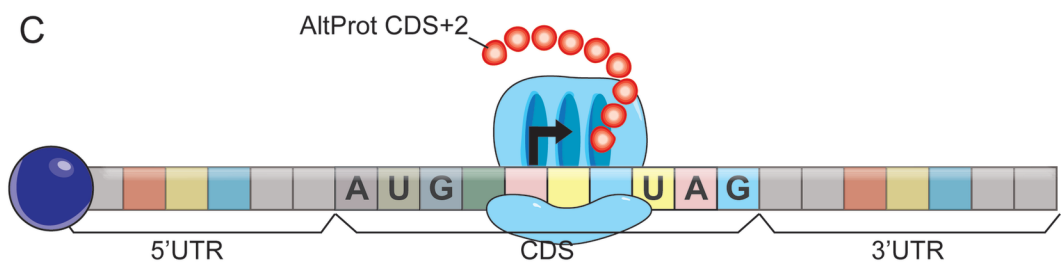
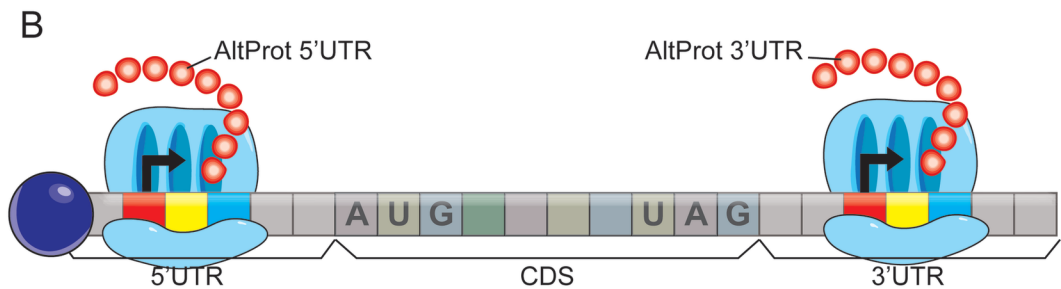
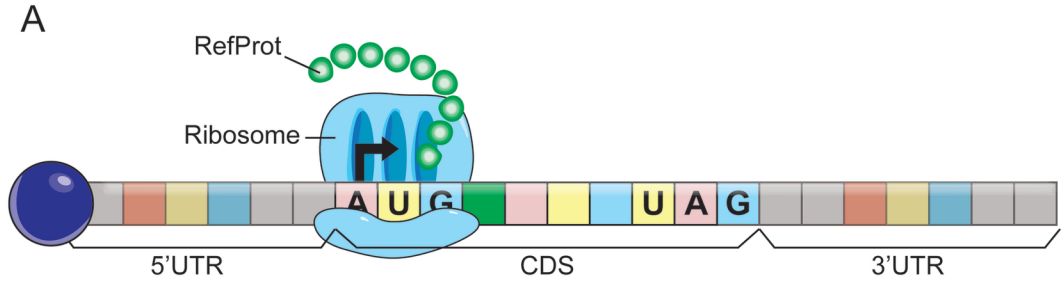
677 **Proteins expression is multi-layered**, involving genes, mRNAs, and allowing the translation
678 of RefProts and AltProts. **A)** The conventional representation admits the expression of a
679 single RefProt derived from an mRNA transcript itself derived from a DNA gene, however
680 the presence of AltProt is now proven. **B)** Genetic mutations impact the transcriptome. If
681 some mutations are silencing mutation to the RefProt they will impact region involved in
682 AltProts translation. The study of AltProt mutations in pathologies is a yet unexplored field
683 that could be highly beneficial and provide massive answers to pathophysiological
684 mechanisms.

685

686 **Figure I.**

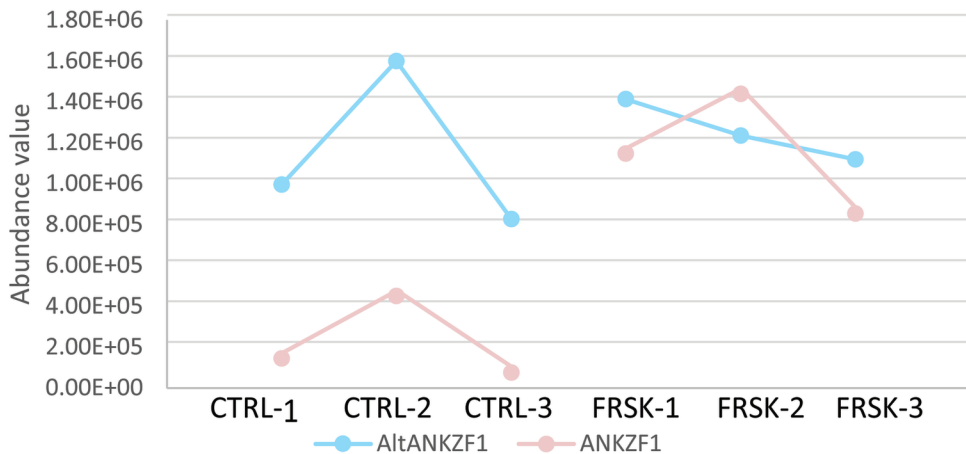
687 **The key events of the Ghost proteome timeline.**

688



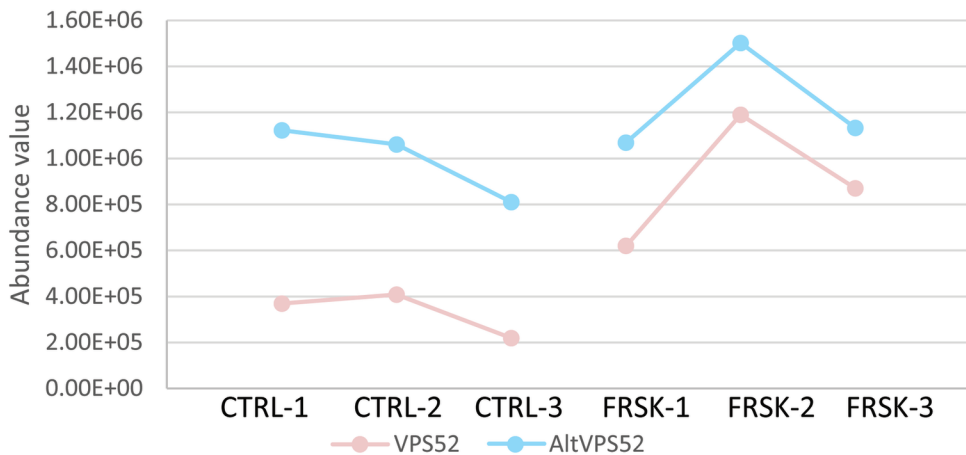
A

Abundance variation of the reference protein:
ANKZF1 and alternative protein: AltANKZF1



B

Abundance variation of the reference protein:
VPS52 and alternative protein: AltVPS52



EMT

Esophageal cancer

AltLINC00624, AltSETD1B
AltMAP2, AltTRNAU1AP
AltEPHA5

YY2BM

**deregulating
cellular
energetics**

**evading growth
suppressors**

HOXB-AS3
Mitoregulin

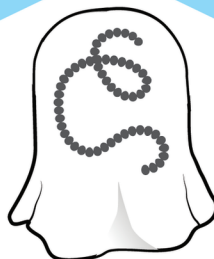
YY2BM

SPAR
Humanin

**resisting
cell death**

**tumor
promoting
inflammation**

Humanin



Ghost Proteome

CASIMO1

**sustaining
proliferative
signals**

**genome instability
& mutation**

MRI-2
AltATF4

ZFAS1

**activating invasion
& metastasis**

CASIMO1

AltCHURC1, AltMDPZ,
AltPIGM, AltSLAMF8,
AltADD3, AltAIG1,

AltMRV11
AltGNL1, AltEDARADD
AltGNL1, AltRP11-576E20.1,

Breast cancer

Ovarian cancer

Glioma cancer

