# Beyond the ribosome: proteome-wide secretability studies with SECRiFY

Boone, M.[1,2,7]; Ramasamy, P.[1,3,4,5,6], Maddelein, D. [1,3], Turan, D. [1,3], Vandermarliere, E.[1,3], Vranken, W.[4,5,6], Callewaert, N.[1,2]

1. Center for Medical Biotechnology, VIB, Zwijnaarde, Belgium

2. L-ProBE, Department of Biochemistry and Microbiology, Faculty of Sciences, UGent, Ghent, Belgium

3. Department of Biochemistry, Faculty of Medicine and Health Sciences, UGent, Ghent, Belgium

4. Structural Biology Brussels, VUB, Brussels, Belgium

5. Structural Biology Research Center, VIB, Brussels, Belgium

6. Interuniversity Institute of Bioinformatics in Brussels (IB)[2], ULB-VUB, Brussels, Belgium

7. Present address: Department of Biochemistry and Biophysics, UCSF, San Francisco, CA, USA

## Abstract

While transcriptome- and proteome-wide technologies to assess processes in protein biogenesis up to the ribosome-associated stages are now widely available, we still lack global approaches to assay post-ribosomal processes, in particular those occurring in the eukaryotic secretory system. We developed SECRiFY to simultaneously assess the secretability of >$10^5$ heterologous protein fragments by two yeast species, *S. cerevisiae* and *P. pastoris,* using custom fragment libraries, surface display and a sequencing-based readout. SECRiFY generates datasets that enable datamining into protein features underlying secretability, and it will enable studies of the impact of secretory system perturbation on the secretable proteome.

## Main text

The eukaryotic secretory system processes roughly a quarter of the proteome, ensuring correct folding, assembly, and delivery of proteins to the extracellular environment, the plasma membrane, or membrane-bound organelles[1–3]. Many proteins and processes involved in these pathways have been identified through studies using just a few model secretory cargos, but an integrative understanding of their role in enabling production of the thousands of secretory proteins is still lacking. For example, it is generally unknown which chaperones are critical for assisting the folding of which types of secretory protein domains. Unfortunately, most current approaches are unsuited to study a comprehensive range of protein folds after entry in the ER. In practice, the absence of such an integrated picture of the secretory system is most apparent in the unpredictability of heterologous protein secretion. Indeed, four decades after the

35  recombinant DNA revolution, obtaining detectable levels of functional protein in a particular
36  heterologous host system (secretory or not) is still principally a process of trial and error. This
37  slows down progress in the many fields of basic and applied life science where recombinant
38  protein production comes into play. Although expression screening of parallel constructs or
39  variant libraries of a protein of interest has gained momentum[4–6] to increase heterologous
40  protein expression success rates, it has to be repeated for each new target. Additionally,
41  alternative comprehensive strategies to assess heterologous expression across entire
42  proteomes have generally been limited to intracellular expression in *E. coli*, small proteomes,
43  or clone-by-clone strategies[7–10].

44  To assay the secretory potential of (human) protein fragments in yeast on a proteome-wide
45  scale, we combined yeast surface display of a novel type of protein fragment libraries and deep
46  sequencing, and dubbed the platform SECRiFY (SECretability screening of Recombinant
47  Fragments in Yeast) (**Fig. 1a** and **e**). By focusing on domain-sized protein fragments rather
48  than full-length proteins, we avoid missing detection of secretability of parts of multi-domain
49  proteins that fail to express or secrete in their entirety due to local issues with misfolding of
50  particular protein areas, translation inhibitory sequences, protease susceptibility, the absence
51  of stabilizing interaction partners or modifications, or toxicity. Chopping up difficult proteins in
52  experimentally tractable fragments has been exploited by structural biologists for years, both
53  in rational target design as well as in random library screens for soluble expression[11–15].
54  Considering the notorious inaccuracies of domain boundary prediction[16], and the knowledge
55  that, even with a reliable estimate, small variations in the exact N- and C-terminus of the
56  fragment can lead to dramatic differences in expressability[16], we opted for a random
57  fragmentation approach. We designed and built directional, randomly fragmented cDNA
58  libraries covering the human transcriptome with fragments coding for approx. 50-100 amino
59  acids, which is the median domain size of human proteins (**Fig. 1b, Suppl. Fig 1a**). To reduce
60  fragment abundance differences inherent to the large dynamic range of mRNA transcripts in
61  human cells, we normalized this library using the Kamchatka crab duplex-specific nuclease[17,18]
62  (**Suppl. Fig. 1b**). Careful optimization of the random primer tag sequence for compatibility with
63  normalization effectively minimized the required library size by approx. 1000-fold, to within a
64  range that is feasibly reconcilable with the diversity that can be routinely obtained in cDNA
65  library construction and yeast transformation (+/- $5*10^6$ - $5*10^7$) (**Fig. 1c-d, Suppl. Results**).
66  To our knowledge, this is the first time that an effective method for normalization of tagged
67  random-primed cDNA fragment libraries is reported, and required substantial development
68  (see also **Suppl. Fig. 1c**). This new method should also find many applications in areas where
69  the protein-coding potential of a cell needs to be effectively covered in expression libraries.

70    Relying on the sophisticated quality control (QC) machinery of the eukaryotic secretory system,

71    which ensures efficient degradation of unstable or misfolded proteins[19,20], we further reasoned

72    that surface display could be used as a proxy for productive secretion. As such, once cloned

73    into the surface display vector and transferred to yeast, library polypeptides are directed to the

74    secretory system by an N-terminal secretory leader sequence (MFα1 prepro), and furthermore

75    on the yeast cell wall via C-terminal fusion to the GPI-anchoring region of the *S. cerevisiae* cell

76    wall protein Sag1 (**Fig. 1a**). Fragments for which the fragment-Sag1 fusion successfully passes

77    (or escapes) secretory system QC without proteolytic degradation are recognized through their

78    N- and C-terminal epitope tags (FLAG and V5, resp.), and are segregated from the rest using

79    iterations of high-efficiency magnetic- and fluorescence-activated cell sorting (MACS/FACS)

80    (**Fig. 1e**). Finally, fragment identification and classification is achieved by deep sequencing of

81    recovered fragment amplicons from the unsorted and sorted cell populations (**Fig. 1e and**

82    **Suppl. Fig 2**).

83    We first benchmarked the method in *S. cerevisiae*. Triplicate secretion screening of a $1.96*10^6$

84    clone fragment library of the HEK293T transcriptome in this yeast revealed that on average

85    $1.76\% \pm 0.12\%$ of library cells displayed a fragment with an intact N-terminus (FLAG-tag) and

86    intact C-terminus (V5-tag) (**Fig. 1e**, **Suppl. Fig 3**). Accounting for a 1/9 chance of up- and

87    downstream in-frame cloning, this means that approximately 15.8% of in-frame fragments are

88    detectably displayed and hence, potentially secretable. After a 32-fold enrichment of these

89    double positive cells through a single round of MACS and two subsequent rounds of FACS

90    (**Fig. 1e**, **Suppl. Fig. 3**), followed by sequencing of the final sorted population, on average,

91    $1.12*10^6$ unique fragments/replicate were detected, covering on average $26.45\% \pm 0.86\%$ of

92    the human canonical transcriptome with at least three reads (**Supp. Table 1-3**). To assess the

93    secretion-predictive value of the method, we picked random clones from the sorted population

94    of a single experiment (**Suppl. Fig. 4, Suppl. Table 4**) and tested the secretion of their

95    encoded fragments when not fused to the anchor protein Sag1. The N- or C-terminal tags of

96    18/20 (90%) fragments could be reliably detected on western blot from the growth medium,

97    and for 16/20 (80%) fragments, both tags were recognized (**Fig. 1f, Supp. Table 5**). As such,

98    fragments displayed by sorted cells are indeed 'secretable' with a high probability. We further

99    classified fragments into those that were enriched (also referred to as secretable) and those

100   that were passively depleted (hence, not detected as secretable) by sorting, setting an arbitrary

101   cut-off on the enrichment factor ($E\ factor = log2\frac{FPTM_{sorted}}{FPTM_{unsorted}}$) at 1 and -1, respectively,

102   reflecting a minimal 2-fold increase and decrease in normalized sequence read counts after

103   sorting. Indeed, as our MACS/FACS based sorting scheme was operated as a binary classifier

104   (sorted/non-sorted), a binary threshold-based classifier is more appropriate in data processing

105   than generalized linear models. Of 170,226 in-frame fragments commonly detected in the three

106  experiments, 6.83% were consistently enriched in all three replicates, and 80.21% consistently

107  depleted (**Supp. Table 6, Suppl. Fig. 5**). Thus, using this metric, these screens were

108  reproducible with an 87.03% concordance between replicates. Even this partial-coverage

109  dataset now represents by far the largest resource on eukaryotic secretability of human protein

110  fragments. To enable researchers to easily access and analyze the data from the screens in

111  this study, we built an interactive database (http://iomics.ugent.be/secrify/search) allowing

112  visualization of the protein fragments detected in these screens and their mapping on available

113  PDB structures (**Suppl. Fig. 11,** Figshare links to data in **Online Methods**).

114  Just as cytosolic protein expression is influenced by a variety of DNA, mRNA, and protein

115  sequence or structural features and their complex interplays[21–24], fragment secretion will

116  depend on a combination of multiple parameters, some of which are related to the unique

117  environment and QC machinery of the ER and beyond. Even already at the simple level of

118  general averaged parameters over the entire secretable vs. non-secretable protein fragment

119  collections, several intriguing observations emerged from our data. Based on predictions of

120  biophysical propensies and on PDB structure mapping (see also **Suppl. Fig. 12**), secretable

121  fragments not only tend to have a lower $\alpha$-helical content than depleted fragments (**Fig. 2a-b,**

122  **Suppl. Fig. 13-14**), but are also distinctly more flexible and intrinsically disordered (**Fig. 2c-d,**

123  **Suppl. Fig. 15a-16**), as well as slightly compositionally biased (**Suppl. Fig. 17-18**, see also

124  **Suppl. Results**). Possibly, this reflects how unstructured sequences that lack typical exposed

125  hydrophobic amino acids are missed by ER chaperones and can subsequently travel

126  downstream. In contrast, no clear differences in number of Cys or N-glycosylation sequons

127  were observed (**Suppl. Fig. 17f-g, Suppl. Fig. 18f-g**). Increasing the fidelity of these findings,

128  all of these patterns were reproduced in screens of a different library of slightly larger fragments

129  in a different yeast species (*P. pastoris*) (**Suppl. Results, Suppl. Table 7-10, Suppl. Figs. 6-**

130  **19**). As for the early folding probability (**Suppl. Fig. 14b**), in cases where depleted and enriched

131  fragments overlap in sequence on the same protein, the presence or absence of regions that

132  are most likely to fold rapidly often correlates with secretability (**Fig. 2e**). Of those fragments

133  that mapped to known structures (roughly 50% of representative fragments), secretable

134  fragments are enriched in distorted sandwich and $\beta$ complex folds compared to depleted

135  fragments, suggesting that these folds are potentially more stable in the secretory

136  environment, while the opposite is true for proteins with, for example, an $\alpha$ horseshoe

137  architecture (**Suppl. Table 11**, Figshare links to data in **Online Methods**). Similarly, certain

138  Pfam domains, such as the AAA18-domain (PF13238), are more prominent in enriched

139  fragments than depleted fragments, while many typically cytoplasmic domains such as

140  ribosomal protein domains or the tetratricopeptide repeat (TPR, PF13181) are found

141  exclusively in depleted fragments (Figshare links to data in **Online Methods,** see also **Suppl.**

142 **Fig. 19**). This illustrates that sequence- and fold-contextual patterns of features still contain
143 much information that is not apparent from averaged parameters. For the first time, our
144 SECRiFY method generates secretability data at a scale at which training of predictive
145 machine learning classifiers becomes feasible and this will be a subject of future work.

146 With SECRiFY, we thus demonstrate that the secretability of protein fragments across entire
147 proteomes can be verified experimentally in an efficient, high-throughput and reproducible
148 manner. Already, the databases we have generated in this work constitute by far the largest
149 resources of such yeast-secretable human protein segments, which will be useful in structural
150 studies (where high levels of proteins are required for crystallization), immunological
151 experiments (where recombinant production is needed for vaccine development or antigen
152 discovery), and biochemical characterizations of particular proteins (which also necessitates a
153 minimal protein amount). Furthermore, it is likely that SECRiFY will provide a means to
154 characterize the substrate scope of secretory system processes that regulate secretory protein
155 passage through the eukaryotic secretory system in a proteome-wide manner. This
156 complements existing methods, such as ribosome profiling[25], which deal with protein
157 biogenesis prior to passage through the secretory system.

158

170

## Author Contributions

172 M.B. performed library constructions, methods development and optimization, screenings,
173 sequencing data processing and analysis, and wrote the manuscript. P.R. performed structure-

174 based data-interpretation under the supervision of E.V. and W.V. W.V. ran, analyzed and

175 visualized the biophysical predictions. D.M. and D.T. developed the website interface and the

176 underlying database. N.C. conceived the project, and assisted in experimental design,

177 interpretation, and manuscript writing. All authors read, revised and approved the final version

178 of the manuscript.

179

## Competing Interests

181 The authors declare that no competing financial interests exist.

182

183
## References

185

186 1.  Barlowe, C. K. & Miller, E. A. Secretory protein biogenesis and traffic in the early secretory
187 pathway. *Genetics* **193,** 383–410 (2013).

188 2.  Braakman, I. & Hebert, D. N. Protein folding in the endoplasmic reticulum. *Cold Spring Harb.*
189 *Perspect. Biol.* **5,** (2013).

190 3.  Aviram, N. & Schuldiner, M. Embracing the void-how much do we really know about targeting
191 and translocation to the endoplasmic reticulum? *Curr. Opin. Cell Biol.* **29C,** 8–17 (2014).

192 4.  Cornvik, T. *et al.* Colony filtration blot: a new screening method for soluble protein expression in
193 Escherichia coli. *Nat. Methods* **2,** 507–509 (2005).

194 5.  Seitz, T. *et al.* Enhancing the stability and solubility of the glucocorticoid receptor ligand-binding
195 domain by high-throughput library screening. *J. Mol. Biol.* **403,** 562–577 (2010).

196 6.  Lockard, M. A. *et al.* A high-throughput immobilized bead screen for stable proteins and multi-
197 protein complexes. *Protein Eng. Des. Sel. PEDS* **24,** 565–578 (2011).

198 7.  Martinez Molina, D. *et al.* Engineering membrane protein overproduction in Escherichia coli.
199 *Protein Sci.* **17,** 673–680 (2008).

200 8.  Luan, C.-H. *et al.* High-throughput expression of C. elegans proteins. *Genome Res.* **14,** 2102–
201 2110 (2004).

202 9.  D'Angelo, S. *et al.* Filtering 'genic' open reading frames from genomic DNA samples for
203 advanced annotation. *BMC Genomics* **12 Suppl 1,** S5 (2011).

204 10. Gupta, A. *et al.* A novel helper phage enabling construction of genome-scale ORF-enriched
205 phage display libraries. *PloS One* **8,** e75212 (2013).

206 11. Reich, S. *et al.* Combinatorial Domain Hunting: An effective approach for the identification of
207 soluble protein domains adaptable to high-throughput applications. *Protein Sci.* **15,** 2356–2365
208 (2006).

209  12.  Yumerefendi, H., Tarendeau, F., Mas, P. J. & Hart, D. J. ESPRIT: an automated, library-based
210       method for mapping and soluble expression of protein domains from challenging targets. *J.
211       Struct. Biol.* **172,** 66–74 (2010).

212  13.  An, Y., Yumerefendi, H., Mas, P. J., Chesneau, A. & Hart, D. J. ORF-selector ESPRIT: a
213       second generation library screen for soluble protein expression employing precise open reading
214       frame selection. *J. Struct. Biol.* **175,** 189–197 (2011).

215  14.  Pedelacq, J.-D. *et al.* Experimental mapping of soluble protein domains using a hierarchical
216       approach. *Nucleic Acids Res.* **39,** e125 (2011).

217  15.  Hart, D. J. & Waldo, G. S. Library methods for structural biology of challenging proteins and
218       their complexes. *Curr. Opin. Struct. Biol.* **23,** 403–408 (2013).

219  16.  Prodromou, C., Savva, R. & Driscoll, P. C. DNA fragmentation-based combinatorial approaches
220       to soluble protein expression Part I. Generating DNA fragment libraries. *Drug Discov. Today* **12**,
221       931–938 (2007).

222  17.  Zhulidov, P. A. *et al.* Simple cDNA normalization using kamchatka crab duplex-specific
223       nuclease. *Nucleic Acids Res.* **32,** e37 (2004).

224  18.  Bogdanov, E. A. *et al.* Normalizing cDNA libraries. *Curr. Protoc. Mol. Biol.* **Chapter 5,** Unit
225       5.12.1-27 (2010).

226  19.  Thibault, G. & Ng, D. T. W. The endoplasmic reticulum-associated degradation pathways of
227       budding yeast. *Cold Spring Harb. Perspect. Biol.* **4,** (2012).

228  20.  Xu, C. & Ng, D. T. W. Glycosylation-directed quality control of protein folding. *Nat. Rev. Mol.
229       Cell Biol.* **16**, 742-752 (2015).

230  21.  Kliman, R. M., Irving, N. & Santiago, M. Selection conflicts, gene expression, and codon usage
231       trends in yeast. *J. Mol. Evol.* **57,** 98–109 (2003).

232  22.  Ghaemmaghami, S. *et al.* Global analysis of protein expression in yeast. *Nature* **425,** 737–741
233       (2003).

234  23.  Zur, H. & Tuller, T. Strong association between mRNA folding strength and protein abundance
235       in S. cerevisiae. *EMBO Rep.* **13,** 272–277 (2012).

236  24.  Boël, G. *et al.* Codon influence on protein expression in E. coli correlates with mRNA levels.
237       *Nature* **529,** 358–363 (2016).

238  25.  Ingolia, N. T., Ghaemmaghami, S., Newman, J. R., S. & Weissman, J. S. Genome-wide
239       analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science* **324**,
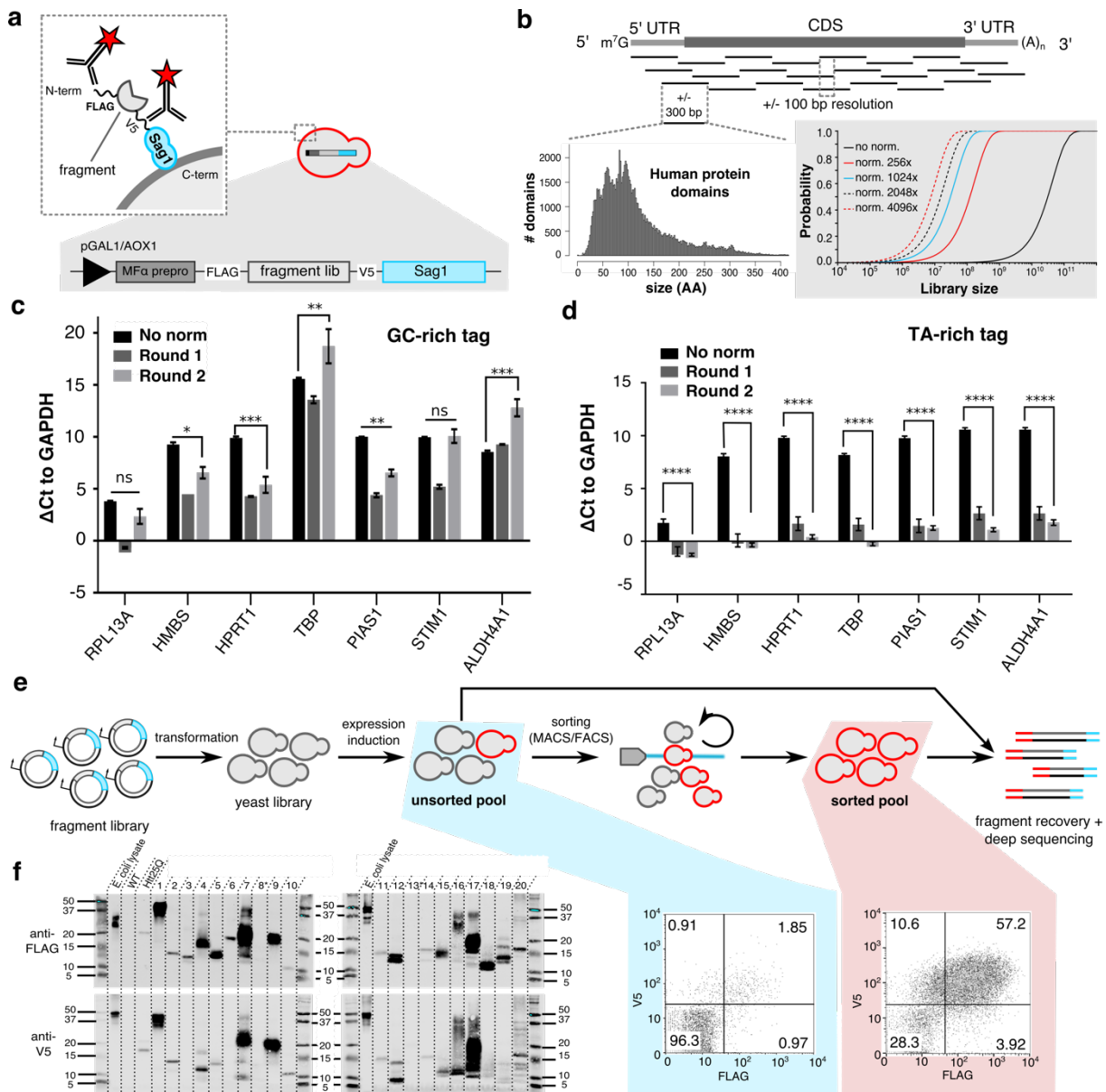240       218–223 (2009).

# Figures



**Fig. 1. The SECRiFY surface display platform for secretability screening of fragment libraries in yeast.** (**a**) Surface display as a proxy for secretion. Productive passage through the yeast secretory system leads to antibody-based labeling of displaying clones through FLAG and V5 epitope tags. (**b**) Human mRNA is fragmented to mimic the size distribution of human protein domains (lower left, Gene3D, n= 104,734). Lower right in grey: estimated relationship between library size and the probability of sampling any fragment, depending on the efficiency of fragment abundance normalization. (**c** and **d**) Normalization of abundance differences compared to *GAPDH*, as ΔCt +/- SEM, is effective only when using a TA-rich (d), not a GC-rich (c) tag in the random primer. Two-way ANOVA with Tukey post-hoc, ns: non-significant, * p<0.05, ** p <0.01, *** p<0.001, **** p<0.0001. n= 9. (**e**) After fragment expression

8

induction, displaying FLAG$^+$V5$^+$ clones are sorted in multiple rounds of MACS/FACS. Fragments are identified by deep sequencing of both pools. (**f**) The majority of protein fragments from sorted cells are detected as secreted when expressed without the Sag1 display anchor.
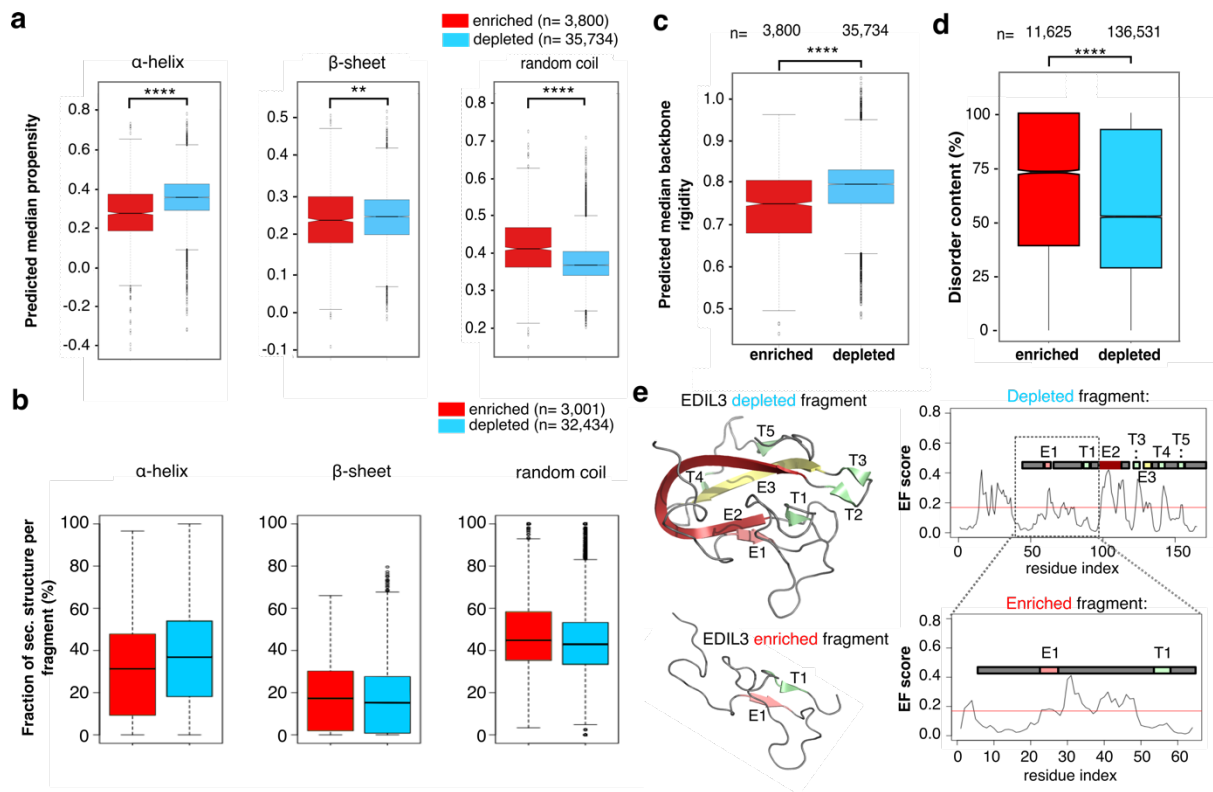
**Fig. 2. Patterns in secretable fragments. (a).** Predictions of secondary structure propensity in enriched and depleted fragments in *S. cerevisiae* shows that enriched fragments have a lower helical content and a higher random coil propensity, which is confirmed further by mapping fragments to known structures in PDB (**b**). Enriched fragments are also predicted to be more dynamic (**c**) and disordered (**d**) than depleted fragments. (**e**). Two overlapping fragments of the human protein EDIL3 differ in secretability outcome. Early folding (EF) propensity predictions suggest that for the depleted fragment regions E2, T3/E3 and R4 are likely the regions driving folding of the depleted fragment, and lack of these regions in the enriched fragment result in a change in secretability.