

1

2

3 **Smoothing membrane protein structure determination by initial** 4 **upstream stage improvements**

5

6 Augusto Quaresma Pedro^{1,2}, João António Queiroz¹, Luís António Passarinha^{1,3,*}

7

8 ¹CICS-UBI – Centro de Investigação em Ciências da Saúde, Universidade da Beira Interior,
9 6201-001 Covilhã, Portugal.10 ²CICECO - Aveiro Institute of Materials, Department of Chemistry, Universidade de Aveiro,
11 3810-193 Aveiro, Portugal.12 ³UCIBIO@REQUIMTE, Departamento de Química, Faculdade de Ciências e Tecnologia,
13 Universidade Nova de Lisboa, 2829-516 Caparica, Portugal.

14

15 *Corresponding author – E-mail: lpassarinha@fcsaude.ubi.pt - Phone: +351 275 329 069;
16 Health Sciences Research Centre (CICS-UBI); Faculty of Health Sciences, University of Beira
17 Interior, Av. Infante D. Henrique, 6200-506, Covilhã, Portugal

18

19

20 **Keywords:** Production; Membrane Protein; Host; Codon usage; Protein 3D-structure; Structure
21 determination; Quality control; Optimization.

22

23 **Abstract**

24 Membrane proteins (MP) constitute 20-30 % of all proteins encoded by the genome of
25 various organisms and perform a wide range of essential biological functions. However, despite
26 they represent the largest class of protein drug targets, a relatively small number high-resolution
27 3D structures have been obtained yet. Membrane protein biogenesis is more complex than that
28 of the soluble proteins and its recombinant biosynthesis has been a major drawback, thus
29 delaying their further structural characterization. Indeed, the major limitation in structure
30 determination of MP is the low yield achieved in recombinant expression, usually coupled to
31 low functionality, pinpointing the optimization target in recombinant MP research. Recently, the
32 growing attention that have been dedicated to the upstream stage of MP bioprocesses allowed
33 great advances, permitting the evolution of the number of MP solved structures. In this review,
34 we analyse and discuss effective solutions and technical advances at the level of the upstream
35 stage using prokaryotic and eukaryotic organisms foreseeing an increase in expression yields of
36 correctly folded MP and that may facilitate the determination of their three-dimensional
37 structure. A section on techniques used to protein quality control and further structure
38 determination of MP is also included. Lastly, a critical assessment of major factors contributing
39 for a good decision-making process related to the upstream stage of MP is presented.

40

41

42 **1. Recombinant membrane protein biosynthesis**

43 Membrane proteins (MP) constitute 20-30 % of all proteins encoded by the genome of
44 various organisms (Lantez et al 2015) and perform a wide range of essential biological
45 functions, thus representing the largest class of protein drug targets (Bernaudat et al 2011).
46 However and despite their biological relevance, most of these targets still do not have any
47 assigned function (Bernaudat et al 2011), as reflected by the relatively low number of MP
48 structures recorded in Stephen White's laboratory database
49 (<http://blanco.biomol.uci.edu/mpstruc/>) - 876 unique MP structures in March 2019. Indeed,
50 determining the structure of a MP is quite complex, mostly due to problems arising from MP
51 low natural abundance, their toxicity when overexpressed in heterologous systems, and
52 difficulties in purifying stable functional proteins and obtaining well-diffracting crystals (Gul et
53 al 2014; Lantez et al 2015). To cope with MP low natural abundance that limits subsequent
54 structural and functional studies, four different approaches have been proposed (Popot 2018),
55 namely: 1) overexpression *in vivo* and *in situ*; 2) overexpression *in vivo* in inclusion bodies; 3)
56 cell-free expression (CFE) *in vitro*; 4) chemical synthesis for short MP or MP fragments. Here,
57 we will generally address the first two approaches based on the following host cells:
58 *Escherichia coli* (*E. coli*), *Pichia pastoris* (*P. pastoris*), also known as *Komagataella phaffii*,
59 mammalian cell lines. The process to obtain a recombinant protein involves the synergy of three
60 key elements – a gene, a vector and an expression host – (Bernaudat et al 2011) and, at least at
61 the theoretical level, is straightforward (Rosano et al 2014). In practice, many things can go
62 wrong, and distinct problems can be found including poor growth of the host, inclusion body
63 formation or lack of protein biological activity (Rosano et al 2014). Indeed, targeting an
64 overexpressed MP to a membrane in such a way they can insert and achieve its native structure
65 is far from being an easy task, once they tend to be toxic, leading to low expression yields of
66 often misfolded and aggregated MP (Popot 2018; Rajesh et al 2011). Moreover, the high
67 diversity of structures and physico-chemical properties displayed by MP makes unfeasible to
68 accurately predict if a protein of interest will express well, be easy to purify, be biologically
69 active or crystallize in any given experimental protocol (Bernaudat et al 2011). Based on the
70 exposed, the development of improved strategies in the recombinant MP production pipeline
71 foreseeing to increase their expression yields in a correctly folded form is crucial in MP
72 research. The evaluation of purified protein quality is crucial in any protein production process
73 and should be accurately performed to avoid irreproducible and misleading observations in the
74 subsequent studies (Raynal et al 2014). After production, MP need to be efficiently solubilized
75 (recently reviewed by Hardy et al 2018 and Popot 2018) and purified (Pandey et al 2016), from
76 which their quality in terms of purity, homogeneity, activity and structural conformity should be
77 assessed (Oliveira and Domingues et al 2018; Raynal et al 2014). In this review, generic

78 guidelines and host characteristics aiming an accurate choice of the host expression system that
79 better suits particular needs will be initially overviewed in this review, and then we discuss
80 important advances reported at the level of the upstream stage of recombinant MP production
81 processes using *E. coli*, *P. pastoris* and mammalian cell lines, representative of major
82 expression systems used for protein expression. Subsequently, general techniques to perform the
83 quality control of the target protein are presented and at the end, insights and directions for a
84 successful MP production pipeline are shown.

85

86 **1.1. Economics vs complexity: guidelines to choose the right host**

87 The most common systems for MP overexpression are microbial (bacteria or yeasts) or
88 higher eukaryotes (insect or mammalian cells) [reviewed in (Bernadaut et al 2011; Fernández
89 and Vega 2016; He et al 2014; Midgett et al 2007; Wagner et al 2006)]. There is no such a
90 perfect host that suits all MP expression projects once they all have advantages and limitations,
91 as highlighted in Table 1. Moreover, the reasons why some MP are overexpressed but others are
92 expressed at low levels are not fully known, although it can be related to how difficult is to fold
93 MP into a functional state (Andréll and Tate 2013).

94 In terms of increasing complexity, the expression systems can be grouped as follows:
95 bacteria < yeasts < insect cells < mammalian cell lines. With an increasing complexity, there is
96 generally an increase in the ability of the host cell to perform native post-translational
97 modifications (PTM). As such, heavily glycosylated proteins are expected to be produced in a
98 more native and folded form from mammalian cell lines, and those obtained from yeasts may
99 not present the native glycosylation profile. On the other hand, simpler hosts such as bacteria
100 allow high productivities, and combine the speed with easiness of operation at a lower cost.

101 Requirements in terms of specific PTM or a near-native-like environment for some
102 mammalian MP are usually the factors dictating the choice of mammalian cell lines, which
103 usually makes use of human embryo kidney (HEK) and Chinese hamster ovary (CHO) cell lines
104 and both cell lines can be applied in stable and transient transfections (He et al 2014; Lyons et al
105 2016). The process of recombinant protein production by transient expression involves the
106 generation of plasmid, transfection in log phase, optional feeds from 24h onwards and then
107 harvest from 48h to 14 days, depending on the target protein, cell line and culture conditions
108 applied (McKenzie et al 2018). Contrasting with transient expression, stably transfected cell
109 lines takes more time (months) and usually requires the stable integration of the recombinant
110 DNA into the host cell genome. Since the expression vector has a gene conferring resistance to
111 an antibiotic, stable integrants can be identified by antibiotic selection; moreover, the integration
112 of the gene into the host genome may be random or the host cell can be engineered to contain a
113 specific sequence recognized by a recombinase that allows targeted integration. Selection of

114 clonal cells is additionally required to identify highly expressing cell lines that are stable under
115 prolonged culture (Andréll and Tate 2013). Transient transfection is quick but after scaling-up,
116 batch-to-batch variability in the amount of protein expressed is often observed; on the other
117 hand, although stable gene expression is initially slower and more technically challenging once
118 a clonal cell line is generated, long-term overexpression can be much more consistent, and the
119 purification of large quantities of supercoiled plasmid DNA for transient-expression is not
120 required (Chaudhary et al 2012; Andréll and Tate 2013). Despite the slow growth rate and
121 usually higher cost, the number of MP structures generated based on such systems has
122 considerably increased, being foreseeable that with the increasing use of cryo electron
123 microscopy for structure determination wherein lesser amount of sample is required (e.g. in
124 comparison with crystallographic studies), mammalian systems will be more frequently used
125 (Lyons et al 2016).

126 Other interesting features to be considered when selecting a host: 1) native intracellular
127 localization of the target protein; proteins that function in specific eukaryotic organelles such as
128 mitochondria, chloroplasts and peroxisomes will generally benefit from expression hosts that
129 possess such organelles (Fernández and Vega 2016); 2) types of lipids of host membranes;
130 hydrophobic mismatch may occur due to differences in lipid bilayer composition and thickness
131 between hosts, as highlighted for the overexpression of eukaryotic MP in bacteria, where the
132 absence of sterols, sphingolipids and poly-unsaturated fatty acids in *E. coli* bilayers poses
133 additional challenges to their proper folding (Snijder and Hakulinen 2016); 3) Construct size;
134 proteins larger than 120 kDa are difficult to be efficiently expressed in *E. coli*, and are typically
135 obtained in very low yields, as inclusion bodies or proteolytically degraded (Fernández and
136 Vega 2016).

137 To overcome the limitations displayed by these *in vivo* expression systems – toxicity,
138 limited membrane space for MP functional folding and inefficient transport and membrane
139 insertion mechanisms -, CFE systems have been reported, which rely on the use of prokaryotic
140 and eukaryotic protein synthesis machinery and related elements to direct protein synthesis from
141 added DNA or mRNA templates (He et al 2011; Henrich et al 2015; Zheng et al 2014). In a
142 different way, the preparation of highly hydrophobic peptides representing functional parts of
143 MPs foreseeing their application onto structural and functional studies can be attained via
144 chemical synthesis (Baumruck et al, 2018). Previously, Fernández and Vega (2016) reported
145 some recommendations on which expression host use for a particular protein.

146

147 **1.2. Upstream strategies to improve membrane protein expression levels and/or**
148 **folding**

149 Membrane protein research strongly rely on recombinant production, which is vital for
150 obtaining high quantities of properly folded proteins for further biophysical and functional
151 testing. While it is difficult to define a set of guidelines generally applicable to all MP., here we
152 review distinct strategies (according Figure 1) that have been used to increase MP expression
153 and/or folding using *E. coli*, *P. pastoris* and mammalian cell-based systems (Summarized in
154 Tables S1, S2 and S3 in Electronic Supplementary Information).

155 **1.2.1. *Escherichia coli***

156 *Escherichia coli* expression systems have been largely investigated for recombinant
157 protein production processes, although with a lower success rate for membrane proteins than for
158 soluble proteins. Aiming to reverse this trend, researchers have driven their efforts to develop
159 enhanced upstream stages encompassing optimizations at the genetic-level, strain engineering or
160 culture conditions, which are reviewed in Table S1 (Electronic Supplementary Information).

161 **1.2.1.1. Genetic-level strategies**

162 The expression of proteins outside their original context can pose additional constraints
163 since they might contain codons that are rarely used in the desired host, come from organisms
164 that use non-canonical code or contain expression-limiting regulatory elements within their
165 coding sequence. The genetic code contains 61 nucleotide triplets (codons) to encode 20 amino
166 acids and 3 codons to terminate translation, and such degeneracy enables many alternative
167 nucleic acid sequences to encode the same protein. Moreover, the frequencies with which
168 different codons are used vary significantly between different organisms, between proteins
169 expressed at high or low levels within the same organism, and sometimes even within the same
170 operon (Gustafsson et al 2004). Indeed, each organism seems to prefer a different set of codons
171 over others, a phenomenon termed as codon bias (Quax et al 2015). Based on these
172 observations, metrics for the frequency of optimal codons were proposed, such as the
173 commonly used codon adaptation index (CAI). The CAI for a certain organism is based on the
174 codon usage frequency in a reference set of highly expressed genes, such as the ones encoding
175 ribosomal proteins and the CAI for a specific gene can be determined by comparing its codon
176 usage frequency with this reference set (Sharp and Li 1986; Quax et al 2015).

177 Different codon biases are also correlated with the amount of the corresponding tRNAs,
178 which vary between organisms; for example, eukaryotes commonly use the AGG codon for
179 arginine, although it is rarely used in *E. coli* (Gustafsson et al 2004). If this exerts a negative
180 effect on heterologous gene expression, then the use of the use of *E. coli* strains over-expressing
181 rare tRNAs (which are commercially available) can improve the yields of target proteins, as
182 previously shown for different constructs of connexin carboxyl-terminal domains attached to
183 their 4th transmembrane domain (Kopanic et al 2013).

184 Moreover, the more codons that a gene contains that are rarely used in the expression
185 host, the less likely it is that the heterologous protein will be expressed at reasonable levels and
186 low levels are exacerbated if the rare codons appear in clusters or in the N-terminal. A strategy
187 to overcome this problem involves sequence re-design by changing the rare codons to codons
188 that more closely reflect the codon usage of the host without modifying the amino acid sequence
189 of the encoded protein (Gustafsson et al 2004). Automated codon optimization algorithms have
190 been developed to design coding sequences optimized for increased expression in certain hosts
191 and codon optimization services are currently offered by DNA synthesis companies, which
192 often rely on confidential algorithms. These algorithms optimize codon usage by maximizing a
193 gene's CAI to match that of the expression host, along with optimizing for some sequence
194 features such as GC content and avoidance of repeats and motifs such as ribonuclease
195 recognition sites, transcriptional terminator sites, Shine-Dalgarno-like sequences, and sequences
196 that lead to strong mRNA secondary structures (Quax et al 2015). On-line tools to gene design
197 such as the OPTIMIZER (<http://genomes.urv.es/OPTIMIZER/>) (Puigbo et al 2007) or to
198 analyze codon usage including the CAIcal (<http://genomes.urv.cat/CAIcal/>) (Puigbo et al 2008)
199 are currently available, among many others which make use of distinct optimization parameters
200 (reviewed by Angov 2011; Gould et al 2014; Parret et al 2016).

201 Based on the rationale that changes in protein structure and function can occur after
202 synonymous codon replacement and that protein structure is DNA sequence-dependent,
203 alternative approaches for synonymous codon design such as the “codon harmonization
204 algorithm” have been proposed, which adapts the codons in a way that the original codon
205 landscape of the gene in the original host is maintained in the expression hosts (Angov et al
206 2008; Quax et al 2015). The authors considered that protein synthesis and folding in *E. coli* is
207 co-translational and that nucleotide sequence-dependent modulation of translational kinetics
208 might influence nascent polypeptide folding. Therefore, in this approach, synonymous codons
209 from *E. coli* were selected that match as closely as possible the codon usage frequency used in
210 the native gene, unless empirical structure calculations show that the codons are associated with
211 putative link/end segments which therefore should be translated slowly (Angov et al 2008).
212 Claassens et al (2017) studied the performance of this codon harmonization algorithm and
213 compared with the wild-type variant and optimized gene variants (resorting to proprietary
214 GeneOptimizer algorithm from GeneArt) using different proton-pumping rhodopsins and
215 enzymes from archaea, bacteria and eukarya. Codon harmonization was performed using a
216 codon harmonizer tool (<http://codonharmonizer.systemsbiology.nl>) based on the harmonization
217 algorithm initially proposed by Angov et al (2008), and uses the codon usage frequency tables
218 for the native and expression hosts, based on all codons in the protein-coding genes annotated in
219 NCBI genome assemblies as inputs. The “codon frequency landscapes” were generated and

220 were evaluated quantitatively based upon a proposed Codon Harmonization Index (CHI), in
221 which a value close to 0 indicates a well-harmonized gene; all harmonized variants have a CHI
222 < 0.1 while all codon-optimized and wild-type variants deviate further from the native codon
223 landscape and consequently present CHI higher than that of harmonized variants (> 0.183). It
224 was additionally observed that transcriptional tuning (in this case by changing the concentration
225 of L-rhamnose) generally improves heterologous production of the distinct variants, although
226 the concentration of rhamnose frequently differs among different codon usage variants of the
227 same protein. In general, harmonization is beneficial for increasing membrane-embedded
228 production compared to wild-type variants for some proteins, for which in this study the wild-
229 type CHI score is also highest (as in the case of leptosphaeria rhodopsin, CHI = 0.279).
230 Moreover, when the codon landscape of the wild-type gene in *E. coli* largely deviates from the
231 landscape in the native hosts, harmonization seems to be a promising approach for increasing
232 MP production (Claassens et al 2017). Recent developments point out that irrespective the
233 algorithm used, using a bicistronic design (in comparison with a monocistronic design) does
234 improve protein production in *E. coli* as it may eliminate the translation initiation as the rate-
235 limiting step of the translation process (Nieuwkoop et al 2019). It should also be also remarked
236 the importance of using updated codon usage tables. In this way, Athey et al (2017) reported a
237 database (available at hive.biochemistry.gwu.edu/review/codon) aiming to present and analyse
238 codon usage tables for every organism with public available sequencing data, and which is
239 being routinely updated to keep up with the continuous flow of new data.

240 Instead of whole sequence optimization, synonymous codon substitutions in the region
241 adjacent to the AUG start may lead to significant improvements in expression, thus
242 circumventing the need to consider whole sequence optimization (Nørholm et al 2013). Indeed,
243 codon usage optimization of the N-terminal guarantee an efficient translation start, which have
244 been proved to enhance human tetraspan vesicle protein (TVP) Synaptogyrin 1 expression in *E.*
245 *coli* (Lów et al 2012). Recently, Saladi et al (2018) developed a data-driven statistical predictor
246 named “IMProve”, which combines a set of sequence derived features resulting in an IMProve
247 score. As this value increases, there is also an increase in the probability of success, *i. e.*
248 selecting a MP that expresses in *E. coli*. Currently, the characterization of an integral MP
249 involves the identification and testing of multiple homologs or variants for expression and the
250 predictive power of “IMProve” enables to enrich for positive outcomes by 2-fold by providing a
251 low-barrier-to-entry (Saladi et al 2018).

252 Throughout the years, codon optimizations have been performed on a first screening
253 basis aiming an increase in the yields of properly folded MP, and with much success without
254 noticeable changes in protein structure and function. However, the increasing understanding of
255 the principles of codon bias and mechanisms of translation have been unveiling yet unknown

256 features. In fact, synonymous codons are known to potentially affect protein expression at
257 various levels and increasing evidences have been showing that translation is affected, leading
258 to dramatic alterations in the conformation and processing of some proteins (Mauro 2018).
259 Overall, codon optimization seems appropriate for some applications, e.g. protein evolution and
260 increasing the expression and/or activity of industrial enzymes; however, for recombinant
261 expression of proteins for therapeutics, we should also aim to maintain the conformation and
262 processing of the natural protein sequences (Mauro 2018).

263 In *E. coli* and due to the higher copy number of the target gene usually achieved with
264 plasmid-based systems, recombinant proteins are typically expressed in *E. coli* from medium to
265 high plasmid copy number (PCN) based on a Col1E derived origin of replication, (Baneyx
266 1999). The PCN is correlated with the recombinant gene dosage and can be accurately
267 determined by quantitative Polymerase-Chain Reaction (qPCR) procedures (Lee et al 2006;
268 Martins et al 2015). A recent study by Jensen et al (2017) provided a systematic approach to
269 identify gene disruptions that increase MP expression in *E. coli* and can be used to improve
270 expression of any protein that poses a cellular burden.

271 Based on the combination of some the above-mentioned strategies, namely “codon
272 harmonization”, use of low copy number vectors with moderate strength, suitable leader
273 sequences, and optimization of cell culture conditions, increased targeting to *E. coli* outer
274 membrane of *Chlamydia trachomatis* major outer membrane protein was observed and the
275 formation of inclusion bodies avoided (Wen et al 2016). On the other hand, prokaryotic
276 expression vectors using the rhaB promoter which are almost completely repressed until
277 induced can be suitable for the expression of toxic proteins (Giacalone et al 2006).

278

279 **1.2.1.2. Strain engineering**

280 Remarkable enhancements in MP expression from *E. coli*-based systems have been
281 achieved with engineered strains due to their improved ability to cope with MP-induced
282 toxicity, more efficient chaperone pathways, different substrate uptake rates or reinforced
283 integrity of intracellular structures, e.g. periplasmic space. Earlier observations have shown that
284 protein (including but not limited to MP) overexpression driven by the T7 RNA polymerase in
285 *E. coli* BL21 (DE3) cells can be limited or prevented by cell death (Miroux and Walker 1996).
286 In this regard, by plating *E. coli* BL21(DE3) cells expressing toxic proteins (oxoglutarate-
287 malate carrier protein from mitochondrial membranes and subunit b of bacterial F-ATPase) in
288 agar plates containing IPTG (for a review of these methods, please refer to Schlegel et al 2017),
289 Miroux and Walker (1996) were able to isolate two isolate two survivors, the mutant host
290 strains C41 (DE3) and C43 (DE3), which have become known as the “Walker strains” and
291 widely used for MP overexpression. Latter studies showed that mutations in the *lacUV5*

292 promoter governing expression of T7 RNA polymerase are the key to improved MP
293 overexpression characteristics of C41 (DE3) and C43 (DE3) strains (Wagner et al 2008). The
294 rationale behind the application of BL21 (DE3) for protein production was that T7 RNA
295 polymerase transcribes faster than *E. coli* RNA polymerase and more mRNA results in more
296 overexpressed protein. However, for most MP, strong overexpression leads to the production of
297 more protein than the Sec translocon can process, thus impairing their insertion into the
298 membrane, which thereby highlights the need to tune MP expression aiming to avoid Sec
299 saturation (Wagner et al 2008). Based on these observations, Wagner et al (2008) engineered a
300 new BL21 (DE3) derivative strain designated Lemo21 (DE3) wherein the activity of the T7
301 RNA polymerase can be precisely controlled by its natural inhibitor T7 lysozyme, which
302 plasmid was under the control of the well-titratable rhamnose promoter (Wagner et al 2008;
303 Schlegel et al 2012). The expression of insertase YidC fused to GFP in the cytoplasmic
304 membrane of Lemo21 (DE3) strain was maximal at 1000 μ M rhamnose, and was additionally
305 demonstrated that this strain is compatible with auto-induction media (Schlegel et al 2012).
306 More recently, Baumgarten et al (2017) isolated the mutant56 (DE3) [Mt56 (DE3)] from
307 BL21(DE3) expressing YidC C-terminally fused to GFP, which allows to evaluate if the
308 produced proteins are being targeted to the cytoplasmic membrane. The authors found that this
309 strain produced several MP in higher levels than C41 (DE3), C43 (DE3) or BL21 (DE3), and its
310 improved performance is attributed a mutation in the gene encoding T7 RNA polymerase in
311 position 305 (C:G – A:T transversion), leading to a single amino acid exchange in T7 RNA
312 polymerase (A102D). Rather than lowering T7 RNA polymerase levels [as with C41 (DE3) and
313 C43 (DE3)], the A102D mutation weakens the binding of the T7 RNA polymerase to the T7
314 promoter governing target gene expression (Baumgarten et al 2017).

315 Envisaging an increase in the amount of membrane-embedded and correctly folded
316 mammalian GPCRs (G protein-coupled receptor), Skretas et al (2012) screened libraries of
317 genomic fragments using two different flow cytometric assays, namely by monitoring the
318 binding of a fluorescently labeled ligand to active GPCR and the fluorescence of GPCR-GFP
319 fusions. These screens allowed the isolation of the genes *nagD* (encoding the ribonucleotide
320 phosphatase NagD), *nlpDA* (encoding a C-terminal truncation of the putative outer membrane
321 lipoprotein NlpD) and the three-gene cluster *ptsN-yhbJ-npr* (encoding three proteins of the
322 nitrogen phosphotransferase system) and was additionally proved that their co-expression leads
323 to a marked increase of membrane-integrated and well-folded GPCR and also a prokaryotic MP
324 (Skretas et al 2012). In general, it seems that the enhanced effect is not due to a direct
325 interaction of these genes with the target proteins, but instead by indirect effects, namely
326 induction of stress responses or changes in the composition of the bacterial periplasm (Skretas
327 et al 2012). Foreseeing the identification of genes whose coexpression can suppress MP-induced

328 toxicity, a genome wide screen identified two potent suppressors, namely *djlA* (encoding the
329 membrane-bound DNAk cochaperone Dj1A) and *rraA* (encoding RRaA), an inhibitor of the
330 mRNA-degrading activity of the *E. coli* RNase E (Gialama et al 2017). *E. coli* strains
331 coexpressing *djlA* or *rraA*, referred as SuptoxD and SuptoxR, respectively, strains were found to
332 have a consistent behavior regarding an enhancement production of distinct MP, namely from
333 mammalian and bacterial origin and with different topologies, and perform better than other
334 commercially available strains (Gialama et al 2017).

335 Another method to mitigate the toxic effect of overexpression is “restrained
336 expression”, in which the production of T7 RNA polymerase and the target gene are controlled
337 by distinct promoters, respectively the arabinose promoter and *T7lac* promoter (Narayanan et al
338 2011). Under “restrained expression” conditions, namely addition of minimal quantities of
339 arabinose (0.01 %) to produce low levels of T7 RNA polymerase and omission of IPTG, aiming
340 to explore the occasional derepression occurring at the *lac* operator site of *T7lac* promoter, an
341 increase of 5 to 25-fold in the expression of homologs of cardiac Na⁺/Ca²⁺ exchanger were
342 obtained, in comparison with IPTG-induction. Moreover, improvements were also found per
343 unit of OD600 nm of cells, indicating that “restrained expression” is associated with decreased
344 cellular toxicity. In general, by reducing the frequency of transcription initiation, protein
345 production is slower, which is unlikely to saturate the biogenesis machinery, thereby providing
346 the explanation for the decreased cytoplasmic aggregation and the attendant cytotoxicity when
347 comparing “restrained” and “rapid” (induction with arabinose and IPTG) expression (Narayanan
348 et al 2011). Nannenga and Baneyx (2011) reported the expression of MP in *Δtig* strains
349 [Transcription factor (TF) deficient] which due to TF inactivation, the signal recognition
350 particle (SRP) has unimpeded access to the nascent transmembrane segment, thus resulting in
351 targeting of MP to the inner membrane, while Yidc overproduction promotes MP insertion and
352 folding in the lipid bilayer.

353 A distinct approach aiming an enhancement in the production of soluble integral
354 membrane spanning proteins relied on engineering *E. coli* wild type AF1000 to reduce the
355 growth rate/substrate uptake rate, accomplished by deletions in the phosphoenolpyruvate
356 carbohydrate:phosphotransferase system (PTS), which is responsible for the uptake of various
357 sugars in *E. coli* (Backlund et al 2011). Distinct mutant strains unable to take up glucose were
358 obtained, and characterized as follows: a defective enzyme IIB^{Man}, which unspecifically
359 controls the uptake of mannose but also allows glucose passage (*ptsM*); a defective enzyme
360 IIBC^{Glc} (*ptsG*), specific for glucose uptake, and the double mutant (*ptsG*, *ptsM*). As a result of
361 the removal of *ptsG*, these mutants display a reduced growth rate at high glucose concentrations
362 but they can grow to high cell densities [although more slowly than BL21(DE3)] since they
363 produce no acetic acid. In general, these strains were able to produce some of the MP in study in

364 relatively larger quantities than BL21 (DE3) but whether this enhanced ability is due to the low
365 growth rate or the lack of acetic acid production was not totally clarified (Backlund et al 2011).

366 Finally, based on the previously published protocols used for MP structure
367 determination, Bruno Miroux research group (Hattab et al 2015) revealed the preferences of *E.*
368 *coli* strain-vector combinations for an optimal use of this expression system and successful
369 production of MP. At that time (June 2014), they found that for the determination of 141 unique
370 non *E. coli* MP structures, 163 expression vector/bacterial hosts were applied, from which T7
371 promoter was dominant (63 %), followed by the arabinose, *tac* and T5 promoter based
372 expression systems (17 %, 9% and 7%, respectively). Moreover, within T7 based expression
373 systems, the host BL21 (DE3) was the most employed, followed by the mutants C43 (DE3) and
374 C41 (DE3), accounting with 40, 18 and 16 MP structures, respectively. Overall, this study
375 shows that C41 (DE3) and C43 (DE3) mutants together with the parental host BL21 (DE3) have
376 contributed significantly for the success of bacterial expression systems in structural biology of
377 MP, in which the mutants have been preferably applied for the production of difficult to express
378 MP. Additional remarks show that IPTG concentration and growth temperature are important
379 parameters complementary to the choice of a bacterial host, and that a high copy number vector
380 should be used with C41 (DE3) to take advantage of the strength of the T7 based expression
381 system, whereas for more difficult MP, the mutant C43 (DE3), especially with low copy number
382 plasmids allows to attenuate the transcription of the target gene (Hattab et al 2015).

383

384 **1.2.1.3. Protein fusion methodologies**

385 Aiming to increase MP solubility and folding or to easily track their expression levels,
386 MP have been expressed with distinct fusion partners (tags) such as SUMO (small ubiquitin-
387 related modified, MBP (maltose-binding protein) or GFP (green fluorescent protein),
388 synthesized either as translational (Zuo et al 2005; Liu et al 2012) or transcriptional fusions
389 (Marino et al 2015). In translational fusions, the N-terminal fusion partners are part of the same
390 protein chain of the membrane protein and can be cleaved off after protein production if any
391 proteolytic cleavage site is introduced. On the other hand, transcriptional fusions exploit the
392 presence of an additional RNA sequence upstream of the mRNA sequence of the target MP,
393 leading to a bicistronic mRNA (Marino et al 2017). As a result, the ribosome produces two
394 distinct protein products during translation, thereby eliminating the need to enzymatically
395 remove the fusion protein during purification (Marino et al 2015). As opposed to translational
396 fusions, transcriptional fusions do not lead to a physical linkage of the fusion protein and MP,
397 which eliminates potential interference of the fusion partner in proper folding and functionality
398 of the target protein (Marino et al 2015; Marino et al 2017). Distinct solubility enhancer tags
399 such as SUMO, MBP, TrxA (thioredoxin) or GST (glutathione-S-transferase) with sizes ranging

400 from 7 to 495 amino acids have been reported (Costa et al 2014). Based on the knowledge that
401 ubiquitin exerts chaperoning properties on fused proteins, translational fusions with the
402 ubiquitin-like protein SUMO were successfully explored towards an enhancement of the
403 solubility and biological activity of the severe acute respiratory syndrome coronavirus (SARS-
404 CoV) MP and 5-lipoxygenase-activating protein (FLAP) (Zuo et al 2005). An additional
405 advantage is that SUMO fusion can be cleaved with high specificity by SUMO protease 1 and
406 generates a protein with the native N-terminal (Zuo et al 2005). On the other hand, Liu et al
407 (2012) evaluated different constructs resorting translational fusions of selenoprotein K
408 envisaging its overexpression in *E. coli* better results were achieved with cytoplasmic MBP over
409 periplasmic MBP and SUMO (Liu et al 2012). In addition to the chaperoning properties
410 displayed by MBP and SUMO, these fusion partners also protect the target proteins from
411 degradation by promoting their translocation from the cytosol to the cell membrane (MBP) and
412 nucleus (SUMO) where less protease content exists (Costa et al 2014). Noteworthy, beyond an
413 increase in the target protein solubility – solubility enhancer -, the natural affinity of MBP
414 towards immobilized amylose resins can also be explored as a purification tool; however, this
415 binding is highly dependent on the nature of the target protein as it can block or reduce the
416 amylose interaction (Costa et al 2014). Translational fusions encompassing a solubility
417 enhancer tag – MBP – and an affinity tag – His-tag – to accomplish the dual purpose of
418 increasing the solubility of MP while exploring their high affinity onto specific affinity
419 chromatographic matrices for purification are feasible, as previously reported for selenoprotein
420 K (Liu et al 2012). A distinct strategy envisaging to target proteins to *E. coli* inner membrane
421 reported by Luo et al (2009) is based on the fusion of a novel partner (P8CBD) to prokaryotic
422 and eukaryotic MP. P8CBD was carefully designed and the DNA encoding 58 amino acid
423 residues of *E. coli* Signal peptidase to provide a second transmembrane segment aiming to
424 extend the protein fusion junction into the periplasmic space, which was selected based on its
425 ability to efficiently establish the desired orientation within the inner membrane (Luo et al
426 2009). A chitin binding domain was also engineered to act as an optional affinity tag or
427 detection epitope while at the fusion junction an enterokinase cleavage site and corresponding
428 FLAG epitope were also incorporated. Overall, by making use of the Signal Recognition
429 Particle (SRP) membrane targeting pathway, the expression and membrane translocation of
430 P8CBD fusion proteins is enhanced (Luo et al 2009). The location of translational fusions is an
431 important factor since they can promote different effects when placed at the N-terminus or C-
432 terminus (Costa et al. 2014). This is better exemplified by the attachment of affinity
433 oligohistidine tags to the periplasmic terminus of *E. coli* transporters, which is detrimental for
434 their expression (Rahman et al 2007). A possible explanation for this relies on a possible
435 interference of oligohistidine sequences with the proper translocation of the adjacent segments

436 of the protein across the membrane during biosynthesis once the charge distribution across
437 transmembrane segments is known to have a profound effect on their orientation (Rahman et al
438 2007). The optimum location of the tag is also influenced by the topology of MP. Although Nⁱⁿ-
439 Cⁱⁿ topologies dominate the membrane proteomes of most organisms, one or both termini of a
440 substantial fraction of MP are located on the extracellular or periplasmic side of the membrane,
441 for which tandem Strep-tag II sequences or oligohistidine tags fused to MBP and a signal
442 sequence should be applied (Ma et al 2015).

443 Unlike translational fusions, there is no need to proceed to the enzymatic removal of
444 transcriptional tags once there is no physical linkage between the target MP and the fusion tag
445 (Marino et al 2017). Marino et al (2015) compared the expression of different proteins using
446 translational and transcriptional fusions of genes coding for the fusion proteins Mistic
447 (membrane-integrating sequence for translation of inner membrane proteins from *Bacillus*
448 *subtilis*), SUMO and a shorter version of YBeL respectively, *mstX*, *sumo* and *ybeL*. They
449 created bicistronic mRNA cassettes where the stop codon of the preceding gene (*mstX*, *sumo*,
450 or *ybeL*) overlaps with the start codon of the target protein, thereby mimicking a common
451 genetic organization observed for bacterial operons (Marino et al 2015). They observed an
452 enhanced expression of MP via transcriptional fusions with *mstX* and *ybeL*, and the cause of this
453 effect cannot be attributed to re-initiation of ribosomes, but instead is most likely attributed to
454 the enhanced translation initiation by a more favorable secondary structure in the transcript
455 (Marino et al 2015).

456 Another major breakthrough within this field in many expression systems was made
457 through fusion of fluorescent reporters such as GFP to the target MP (Drew et al 2001;
458 Goehring et al 2014; Gul et al 2014), which behaves as a folding indicator of the target MP and
459 allowing to infer on their expression levels. This process usually relies on fusing GFP to the C-
460 terminal of proteins; since GFP only becomes fluorescent if the MP integrates in the
461 cytoplasmic membrane, it allows to distinguish between MP overexpression in the cytoplasmic
462 membrane and in inclusion bodies at any stage during overexpression, solubilization and
463 purification (Drew et al 2001; Drew et al 2006). In addition, GFP will only become fluorescent
464 if the MP has a Cⁱⁿ topology, i. e. the C-terminus is cytoplasmic (Drew et al 2006). Noteworthy,
465 fluorescence in whole cells can be detected with a detection limit as low as 10 µg of GFP per
466 liter of culture, and can also be determined in standard SDS polyacrylamide gels with a
467 detection limit of less than 5 ng of GFP per protein band (Drew et al 2006). Also based on the
468 use of GFP as a fusion partner, Nji et al (2018) recently reported a fluorescence detection size-
469 exclusion chromatography-based thermostability assay (FSEC-TS) that allows measuring
470 apparent melting temperatures (T_m) of MP in the absence and presence of distinct lipids, which
471 can be helpful to identify which lipids can have a stabilizing effect for a particular target.

472 In addition to GFP, Gul et al (2014) reported the translational fusion of the
473 erythromycin resistance protein (23 S ribosomal RNA adenine N-6 methyltransferase, ErmC)
474 (in tandem with GFP) to the C-terminus of different bacterial MP wherein GFP fluorescence
475 was applied to report the folding state of the target protein and ErmC to select for increased
476 expression. Evolved strains termed NG were selected in increasing concentrations of
477 erythromycin which carry out a mutation in *hns* gene, and the degree of MP expression
478 correlates with the severity of *hns* mutation, although its deletion resulted in an intermediate
479 expression. Overall, in each NG strain, the amount of fluorescent (folded) protein and the ratio
480 of folded over misfolded protein increased up to 10-fold relative to the parental strain
481 BW25113B (Gul et al 2014). Another approach to easily detect the expression levels of MP was
482 reported by Hsu et al (2013) which is based on the use of mutated bacteriorhodopsin from
483 *Haloarcula marismortui* as a fusion partner, and which unlike GFP, MP overexpression can be
484 detected by naked eye or by directly monitoring their optical absorption.

485 Aiming to select mutants of *E. coli* that improve MP expression, Massey-Gendel et al
486 (2009) reported an approach that relies on fusing the targeted MP to a C-terminal selectable
487 marker that confers a drug resistance phenotype (Massey-Gendel et al 2009). The rationale
488 behind this strategy is that the production of the selectable marker and survival on selective
489 media is linked to expression of the targeted MP, namely when the c-terminus is in the
490 cytoplasm. After the selection of the mutants, curing of isolated mutants is performed by *in vivo*
491 digestion with the homing endonuclease I-CreI (Massey-Gendel et al 2009).

492 Recently, Mizrachi et al (2015) developed a technique called SIMPLEx (Solubilization
493 of Integral MP with high Levels of Expression), which allows the direct expression of soluble
494 products in living cells by fusing the target MP with the carboxyl terminal of apolipoprotein A-1
495 (ApoAI*). In addition, a highly soluble “decoy” protein from *Borrelia burgdorferi*, namely the
496 outer surface protein A (MBP lacking its N-terminal signal peptide can also be used) was fused
497 to the N-terminus to prevent the *E. coli* secretory pathway to introduce the protein in inner
498 membrane. Acting as an amphipatic proteic “shield” which sequester MP from water, ApoAI*
499 promotes the solubilization of structurally diverse MP (bitopic α -helical, polytopic α -helical and
500 polytopic β -barrel) and yields of EmrE-solubilized dimers and tetramers (EmrE basic functional
501 units) ranged between 8 and 10 mg/L of culture after Nickel affinity chromatography. ApoAI*
502 solubilized EmrE (*E. coli* ethidium multidrug resistance protein E) was amenable to structural
503 characterization including negative staining electron microscopy, dynamic light scattering and
504 SAXS (Small angle X-ray scattering) data collection (Mizrachi et al 2015).

505

506 **1.2.2. *Pichia pastoris***

507 **1.2.2.1. Genetic-level strategies**

508 Yeasts and particularly *P. pastoris* are highly attractive alternatives for MP expression
509 as they represent low-cost cultivation and high-quantity production platforms, meeting the
510 demand for criteria of safety and authentically process proteins (Emmerstorfer-Augustin et al
511 2019). *Pichia pastoris* systems usually rely on the use of integrative plasmids containing the
512 gene of interest which are integrated into the yeast genome, generating stable production strains
513 (Dilworth et al 2018). Moreover, protein production is usually accomplished resorting the
514 alcohol oxidase promoter (AOX), which is inducible by methanol and depending on the
515 functionality of 1 or both *aox* genes, recombinant strains may present a Mut^S or Mut⁺ phenotype
516 exhibiting different growth behaviors (in methanol) and different methanol requirements for
517 induction. Other commonly used promoter is the constitutive glyceraldehyde-3-phosphate
518 (GAP) dehydrogenase promoter (Gonçalves et al 2013).

519 In the last years, studies have shown that distinct recombinant gene dosages and codon
520 usage optimizations greatly influence MP expression levels in *P. pastoris*. As mentioned above,
521 *P. pastoris* expression systems usually rely on expression plasmids that are integrated into the
522 yeast genome and multi-copy clones – the so-called “jackpot clones” –, can be selected
523 experimentally by screening several colonies in increasing concentrations of antibiotic
524 (Dilworth et al 2018). Nordén et al (2011) performed a two-step antibiotic selection, initially
525 with 100µg/mL zeocin and then with higher concentrations, from which they isolated multi-
526 copy clones and observed that the expression of different aquaporins strongly respond to an
527 increase in recombinant gene dosage, independently of the amount of protein expressed from a
528 single gene copy clone. However, despite higher recombinant gene dosages can lead to higher
529 titers of recombinant proteins, this correlation is not always linear and strains with low copy
530 number may be preferred (Aw and Polizzi 2013, Dilworth et al 2018). Aiming to exclude
531 possible false-positives while establishing accurate correlations, along with the levels of the
532 target protein, the recombinant DNA levels must be evaluated, for which qPCR protocols have
533 been reported using pPICZ vectors (Nordén et al 2011) and resorting to SYBR Green or
534 TaqMan (Abad et al 2010). Another way to improve human aquaporins expression in *P.*
535 *pastoris* is based on the optimization of the nucleotide sequence around the initial ATG based
536 on the use of mammalian Kozak’s sequence consensus (Oberg et al 2009). The prevalence of a
537 guanine at the first position of the second codon after ATG encodes small amino acids such as
538 alanine (GCN) or on a smaller extent glycine (GGN), which are crucial to ensure an efficient
539 cleavage of the initiator methionine (Oberg et al 2009). In most cases, this has a positive impact
540 on aquaporins expression, while the opposite seems to be observed when a thymine is at
541 position +6 (Oberg et al 2009).

542 The codon bias problem in MP production from *P. pastoris* have also been addressed.
543 Considering that the translation efficiency of more highly expressed genes may be especially

544 sensitive to codon usage, Bai et al (2011) generated a codon usage table specific for highly
545 expressed genes in *P. pastoris* and adjusted the sequence of P-glycoprotein-encoding *mdr3*
546 gene, taking into account relative codon frequencies for each amino acid, as well as optimizing
547 GC content and controlling for mRNA instabilities. Using the optimized gene construction, the
548 authors obtained an increase of three-fold in the expression yields in comparison with the wild-
549 type gene of P-glycoprotein and similar secondary and tertiary structures between the proteins
550 from the different constructs, emphasizing the effectiveness of the gene optimization approach
551 developed (Bai et al 2011).

552 Expression resorting fusion partners has been applied since the early beginning of MP
553 expression in *P. pastoris*. Talmont et al (1996) expressed the μ -opioid receptor fused with *S.*
554 *cerevisiae* α -mating factor aiming to facilitate the translocation of the receptor to the membrane.
555 Distinctly, it was shown that the presence of the α -mating factor can be detrimental for the
556 expression of human histamine H₁ receptor in *P. pastoris* (Shiroishi et al 2011), which can be
557 due to incomplete processing by the endogenous Kex2 protease, leading to a heterogenous
558 population. A way to overcome this problem is by introducing a proteolytic cleavage site
559 upstream of the gene (Byrne 2015).

560 The application of GFP as a fusion partner has been extensively used to screen for high-
561 yield expressing clones spanning the most popular hosts for MP production including *P.*
562 *pastoris*. Brooks et al (2013) reported a fluorescent-based induction plate assay aiming the
563 simultaneously screening of *P. pastoris* clones for the expression of aquaporin 4 and
564 homologues of ER associated MP phosphatidylethanolamine N-methyltransferase in which 50
565 and 48 clones were respectively screened. The plates were imaged under blue light and the
566 colony fluorescence quantified using Mean Gray Values and revealed a distribution of
567 fluorescence related to protein expression, ranging from background to high, being additionally
568 demonstrated that there is a good correlation between plate expression and liquid culture
569 expression (Brooks et al 2013).

570 In addition to secreted proteins, MP can also enter the secretory pathway but unlike
571 them, MP remain in the ER, Golgi or the plasma membrane (Vogl et al 2014). Due to MP
572 overexpression, unfolded and misfolded proteins can accumulate in the ER, thereby triggering
573 the unfolded protein response (UPR). The UPR signaling pathway involves the kinase/RNase
574 Ire1 that when activated initiates an unconventional splicing reaction of the *HAC1* mRNA that
575 ends with the removal of the intron and subsequent translocation of Hac1p to the nucleus
576 (Guerfal et al 2010). Guerfal et al (2010) showed for the first time the beneficial effect of co-
577 expressing Hac1p with the adenosine A2A receptor, namely in terms of a better processing of
578 the alpha-mating factor, thus improving the homogeneity of the obtained MP fractions. Later,
579 Vogl et al (2014) performed a transcriptomic analysis of *P. pastoris* CBS 7435 overexpressing

580 different classes of MP (mitochondrial, ER/Golgi and plasma-membrane localized) and found
581 that proteins targeted to the mitochondrial membrane mainly alter the energy metabolism while
582 the gene coding for Hac1p was upregulated in strains expressing the CMP-Sialic acid
583 transporter, which localizes to ER and Golgi. Interestingly, they found that the overexpression
584 of the spliced variant of Hac1 led to an increase of 1.5-fold to 2.1-fold in the expression of ER-
585 resident MP tested (Vogl et al 2014)

586

587 **1.2.2.2. Strain engineering and improved processing conditions**

588 *Pichia pastoris* expression strains are derivatives of NRRL-Y 11430 (Northern
589 Regional Research Laboratories, Peoria, IL, USA) (Cregg et al 2000) encompassing distinct
590 genotypes/phenotypes, and generally most of them have been applied for MP production,
591 namely X33 (wild-type/Mut⁺) (Oberge et al 2009), KM71H (*arg4aox1::ARG4/Mut^S Arg⁺*) (Bai
592 et al 2011), GS115 (*his4/Mut⁺ His⁻*) (Guerfal et al 2010) and also protease deficient strains such
593 as SMD1163 (*pep4 prb1 his4/Mut⁺ His⁻*) (André et al 2006).

594 The requirement of association with cellular membranes and the type of membranous
595 lipids can be critical for successfully achieving the goal of producing a recombinant MP in a
596 functional active form, given their close spatial interactions (Emmerstorfer-Augustin et al
597 2019). Plasma membranes are generally constituted by a mixture of lipids including
598 phosphatidylcholine, phosphatidylethanolamine, phosphatidylinositol, phosphatidylserine,
599 phosphatidic acid, sphingolipids and sterols (van der Rest 1995). As the composition and
600 molecular properties of the lipids differ from lower to higher eukaryotes, the distinct type of
601 sterols in yeasts and mammals, respectively ergosterol and cholesterol, can represent a
602 bottleneck for the heterologous expression of mammalian proteins in yeasts (Emmerstorfer-
603 Augustin et al 2019; Hirz et al 2013). Therefore, aiming an improvement in the functional
604 expression, stability and translocation of Na⁺/K⁺ ATPases $\alpha 3\beta 1$ isoform, Hirz et al (2013)
605 reprogrammed *P. pastoris* (strain CBS7435 $\Delta his4 \Delta ku70$) to mainly produce cholesterol instead
606 of ergosterol. This was accomplished by replacing ERG6 (encodes the sterol C-24 methyl
607 transferase) and ERG5 (encodes the sterol C-22 desaturase) by constitutive DHCR7 and
608 DHCR24 (dehydrocholesterol reductases) overexpression cassettes, envisaging an efficient
609 conversion of cholesta-5,7,24(25)-trienol to cholesterol (Hirz et al 2013; Emmerstorfer-
610 Augustin et al 2019). The authors found that the expression levels of the target ATPase
611 significantly increased with induction time in the cholesterol-forming strain compared to the
612 wild-type strain, indicating a positive influence of the altered sterol composition on the stability
613 of the synthesized MP (Hirz et al 2013). Another example of “humanizing” *P. pastoris* for the
614 expression of human proteins consists of the disruption of an endogenous glycosyltransferase
615 gene (OCH1) and the stepwise introduction of heterologous glycosylation enzymes, envisaging

616 to largely eliminate the fungal N-type N-glycosylation while avoiding a considerable
617 heterogeneity in the produced protein and their rapid clearance if therapeutics is the main goal
618 (Jacobs et al 2009; Laukens et al 2015). This strategy is generally referred as GlycoSwitch® and
619 can be applied in wild-type strains (e.g. GS115) or GlycoSwitch® Man 5 strain wherein the first
620 glyco-engineering step was already introduced, and encompasses distinct glyco-engineering
621 steps based on the transformation of *P. pastoris* with GlycoSwitch® vectors under previously
622 reported protocols (Jacobs et al 2009; Laukens et al 2015). Currently, these vectors are
623 commercially available from BioGrammatics (Carlsbad, USA) under the licence from Research
624 Corporation Technology (RCT).

625 Envisaging to prevent a possible inhibition of the AOX promoter by glycerol, *Pichia*
626 *pastoris* AOX-based bioprocesses usually encompass an initial stage of growth in glycerol
627 followed by methanol induction, which is often cumbersome specially when glycerol
628 consumption cannot be monitored (Lee et al 2017). Earlier observations with KM71H strains
629 demonstrating that leaky expression is not a critical factor once the target expression per cell
630 mass is mostly dependent on the starting glycerol concentration of the media and to a lesser
631 degree by yeast nitrogen base (YNB) and biotin concentrations. Moreover, as even in the
632 presence of a methanol concentration higher than the glycerol concentration no target
633 expression was detected until about 24 h of incubation, Lee et al (2017) developed the Buffered
634 extra-YNB Glycerol Methanol (BYGM) auto-induction media (100 mM potassium phosphate
635 pH 6.0, 2.68 % w/v YNB, 0.4 % v/v glycerol, 0.5 % v/v methanol and 8×10^{-5} % w/v biotin).
636 This auto-induction method avoids the traditional media-swabbing step and it is additionally
637 claimed that it can be applied to Mut^S and Mut⁺ strains and distinct MP without compromising
638 their expression yields (Lee et al 2017). The use of additives in culture media have also been
639 reported to increase MP expression levels. André et al (2006) reported increased expression
640 levels of functional GPCR resorting the optimization of growth temperature and
641 supplementation of culture media with specific GPCR ligands, histidine, and dimethylsulfoxide
642 (DMSO). As DMSO can modify the physical properties of membranes and upregulates genes
643 involved in lipid synthesis (Murata et al 2003), it can have a positive effect on MP in yeast and
644 is additionally pointed out that by permeabilizing membranes, it can have an indirect effect by
645 facilitating the entry of other ligands to intracellular compartments where they reach the
646 receptor populations (André et al 2006). The beneficial effect of DMSO is not restricted to
647 GPCR as Pedro et al (2015) reported an increase of 1.8-fold in the enzymatic activity of human
648 membrane-bound catechol-*O*-methyltransferase (MBCOMT), achieved by adding 5% v/v
649 DMSO. Subsequently, the artificial neural network modelling of the methanol induction phase,
650 accomplished by tailoring the temperature, DMSO concentration and methanol constant flow-
651 rate allowed an improvement of 1.53 fold in the enzyme activity over the best conditions

652 performed in the DoE step (Pedro et al 2015). In addition, the direct solubilization of MP whole
653 cells (yeasts protoplasts) may help to decrease the amount of misfolded and/or aggregated
654 proteins that are co-extracted with the properly folded protein (Hartmann et al 2017).

655

656

657

658 **1.2.3. Mammalian cell lines**

659 General approaches and factors for successful optimization of mammalian-based
660 systems for recombinant protein production have been reviewed elsewhere (Andréll and Tate
661 2013; Almo and Love 2014; Hacker and Balasubramanian 2016; McKenzie 2018). In this sub-
662 section, we will generally focus our attention in strategies that have been proved to be
663 particularly useful for MP, foreseeing improved expression and/or folding and also those
664 enabling biochemical and functional studies of these relevant drug targets (summarized in Table
665 S3).

666 Distinct mammalian cell lines have been applied for MP production such as HEK293,
667 baby hamster kidney cells (BHK-21), monkey kidney fibroblast cells (COS-7) and CHO
668 (Andréll and Tate 2013), but HEK293 and CHO are more commonly applied, either in transient
669 or stable transfection (Lyons et al 2016).

670 The levels of expression of MP in transiently transfected mammalian cell lines are
671 affected by the plasmid size, the amount of plasmid used per transfection, the strength of the
672 promoter, the cell type, the efficiency of the transfection and potentially the toxicity of the
673 transfection reagent (Andréll and Tate 2013). Using design of experiments, Bollin et al (2011)
674 optimized the yields of an antibody resorting to transient gene expression and found that the
675 DNA concentration can be maintained at relatively low concentrations (1 mg/L range). Indeed,
676 envisaging functional expression of a MP in the plasma membrane, the ratio of plasmid DNA
677 added per reaction can be a crucial factor (particularly if a strong promoter is used), once too
678 much plasmid can lead to intracellular accumulation of the protein and potentially misfolded
679 (Andréll and Tate 2013). Both CHO and HEK cell lines have been extensively used in transient
680 transfection, advances in serum free media formulations allow their growth to high-cell
681 densities, which can greatly facilitate the purification of target proteins (Almo and Love 2014;
682 McKenzie et al 2018). An alternative approach increasingly applied as a gene delivery
683 methodology for protein production is based on the use of lentivirus, owing to their ability to
684 transduce a broad range of cell types (Bandaranayake and Almo 2014). Aiming to combine the
685 ease and speed of transient transfection with the robust expression of stable cell lines, Elegheert
686 et al (2018) constructed a lentiviral plasmid suite around the transfer plasmid pHR-CMV-TetO₂
687 that is designed for large-scale protein expression from HEK293 cell lines and allows

688 subcloning of cDNA from the plasmid PHLsec usually applied for transient transfection. This
689 approach was tested in both soluble and MP, and in general, the typical lead time for protein
690 production using this strategy is of 3-4 weeks and approximately three- to tenfold improvement
691 in protein production yield per cell was obtained, in comparison with transient transfection
692 (Elegheert et al 2018).

693 Unlike transient transfection, stable gene expression requires the screening of clonal cell
694 lines, which is typically achieved through limited dilution involving serial dilution of recently
695 transfected cells and seeding on tissue culture plates with antibiotic-resistance media.
696 Subsequently, different colonies are individually transferred to 24-well plates and scaled-up
697 (Andréll and Tate 2013). For a review of selection methodologies, please refer to Browne and
698 Al-Rubeai (2007).

699 Along the years, aiming to easily ascertain the quality and level of expression of target
700 MP, methodologies resorting to GFP fusions have been reported. Particularly, the expression of
701 GFP fused to the termini of MP have been applied to directly monitor in whole cells for their
702 subcellular locations by fluorescence microscopy (Goehring et al 2014). A slightly different
703 approach was reported by Mancina et al (2004), where the production of the target MP and GFP
704 is based on a bicistronic mRNA, thus leading to the production of two separate proteins wherein
705 the high-yielding clones are selected based on a fluorescence-activated cell sorting procedure.

706 Given the relevance of MP as drug targets for a variety of human diseases, advances in
707 mammalian cell-based systems have allowed performing functional studies that otherwise could
708 be highly hampered. Baculovirus mediated gene transduction of mammalian cells (BacMam)
709 has been widely used due to its compatibility with a variety of mammalian cell lines and the
710 possibility of co-infecting with multiple BacMam viruses to express protein complexes (Lyons
711 et al 2014). Shukla et al (2012) exploited this strategy towards the development of a transient
712 expression system for co-expression of two drug transporters (ABCB1 – P-glycoprotein – and
713 ABCG2) in mammalian cells, which is useful to determine their contribution to the transport of
714 a common anticancer drug substrate. Moreover, both transporters were functionally active when
715 co-expressed (Shuka et al 2012). A distinct approach involves the codon-optimization of the
716 sequence of the human sodium/iodide symporter (NIS) based on the highest usage frequencies
717 in humans, while RNA instability motifs, very high (>80%) or very low (<30%) GC content
718 regions and cis-acting motifs were also removed (Kim et al 2015). As a result, the CAI was
719 highly improved (0.79 vs 0.97 for wild-type and optimized sequences) and from transfected
720 cancer cells, it was found that the levels of NIS were enhanced as well as the radioiodine
721 uptake. These results show the importance of codon usage optimizations in the development of
722 more efficient reporters and efficient therapeutic genes, distinct goals than improving MP
723 heterologous expression (Kim et al 2015).

724 To facilitate MP production for structural analysis relies on the use of HEK293S GnTI-
725 (lacking the gene N-acetylglucosaminyl transferase I - GnTI) and a tetracycline-inducible
726 promoter (Chaudhary et al 2012). If on one hand, the lack of GnTI restricts N-linked glycans to
727 a homogeneous Man₅-GlcNac₂, since N-linked glycosylation is often regarded as a barrier
728 toward structure determination via X-ray crystallography due to the heterogeneity and
729 conformational flexibility of these glycans, the inducible promoter allows the establishment of
730 high-density cell cultures which are not always achieved if the target protein tends to be
731 cytotoxic (Chaudhary et al 2012). Alternative approaches have been suggested to overcome
732 toxicity issues associated with MP overexpression. Ohsfelt et al (2012) designed an anti-
733 apoptosis strategy involving co-expression of *Bcl-xL* gene (encodes for an anti-apoptotic
734 protein) aiming to prevent cell death by bioreactor stresses, nutrient depletion, toxin
735 accumulation, and stresses due to folding and processing requirements for complex proteins
736 such as MP. The authors observed that cell death are diminished due to the co-expression of the
737 anti-apoptotic gene and transient production of two different receptors were improved (Ohsfeldt
738 et al 2012).

739
740

1.3. Protein Quality Control:

741 The purity and integrity of purified protein samples are usually evaluated by
742 electrophoresis (native or denaturant) coupled with detection methods with varying sensitivities
743 (Oliveira and Domingues et al 2018; Raynal et al 2014). On the other hand, isoelectric focusing
744 and capillary electrophoresis have also been used to distinguish the protein of interest from
745 closely related undesired subproducts or contaminants (Raynal et al 2014), while UV-Visible
746 spectroscopy is useful to detect nucleic acid contamination (Oliveira and Domingues et al
747 2018).

748 Mass spectrometry (MS) has been widely applied to measure molecular weights of
749 proteins while allowing protein identification by peptide mass fingerprinting (PMF) and based
750 on MS/MS spectra (Zhang et al 2010). By detecting mass changes introduced by post-
751 translational modifications, MS can also be used analyze these modifications (Zhang et al
752 2010). MS-compatible detection methods enable MS analysis after electrophoresis (Raynal et al
753 2014). Despite such analysis are usually performed after purification, Gan et al (2017) reported
754 a native MS approach that allows the characterization of overexpressed recombinant proteins
755 directly in crude *E. coli* lysates, allowing obtaining information on its identity, solubility,
756 oligomeric state, overall structure and stability without purification. Cells were lysed in a buffer
757 supplemented with 1M ammonium acetate to ensure compatibility with MS. Spectra were
758 acquired for distinct proteins with molecular weights ranging from 19 to 47 kDa, and revealed
759 highly resolved peaks, narrow charge state distributions and the anticipated stoichiometry,

760 thereby confirming that at least for these proteins, purification is not a prerequisite (Gan et al
761 2017).

762 In addition to the integrity and purity of the protein sample, homogeneity is also crucial
763 to infer on the correct oligomeric structure of the protein. Dynamic light scattering (DLS) and
764 more accurately analytical size exclusion chromatography (SEC) are useful to these
765 determinations (Oliveira and Domingues et al 2018; Raynal et al 2014). In quality control
766 methodologies, studying the secondary and tertiary structure of proteins is important to infer
767 about their folding and monitor protein conformational changes. A range of spectroscopic
768 techniques have been developed for such task, being circular dichroism particularly useful to
769 determine the secondary structures and folding properties of recombinant proteins (Oliveira and
770 Domingues et al 2018). Based on several generic or protein-specific functional assays which
771 depend upon catalytic and binding properties of the protein of interest, it is also important to
772 determine the activity of the target protein samples (Raynal et al 2014). Additional details of
773 distinct analytical methods used for the characterization of therapeutic proteins including
774 advantages and drawbacks as well as the type of information delivered from each technique can
775 be found in the recent review by Fuh et al (2016).

776

777 **1.4. Insights for better decision-making processes in the upstream stage of** 778 **membrane proteins:**

779 In this review, we addressed the first stage and, more specifically their (bio)synthesis by
780 recombinant production processes. *E. coli*, *P. pastoris* and mammalian cell lines were selected,
781 given their wide applicability and to cover hosts with different inherent complexities. Based on
782 the information here reviewed, general insights to understand which host may better fit in a
783 specific project are presented in the next paragraphs and summarized in Table 2 and Figure 2. *E.*
784 *coli* is probably the better characterized host for which there are many genetic tools available. It
785 is more suitable for low molecular weight MP and is capable to grow easily to high-cell
786 densities at a relatively low cost. Unlike *E. coli*, and mammalian cell lines allow the production
787 of larger MP and protein complexes with proper PTM including glycosylation patterns,
788 although in this regard the performance of mammalian cell lines is best. However, obtaining
789 recombinant proteins which better resemble their native counterparts comes with a cost and
790 these systems are more technically challenging and this process can be lengthy. The
791 methylotrophic *P. pastoris* gathers characteristics from both prokaryotic and the other higher
792 eukaryotic hosts. Particularly, direct and indirect evidences point out the importance of *P.*
793 *pastoris* host membranes wherein the type of lipids can influence the expression yields and
794 overall folding of heterologous human MP while inducing membrane proliferation (HAC1
795 overexpression and possibly the use of DMSO as an additive in culture media). The

796 identification of genes limiting MP overexpression resorting systems biology approaches based
797 on -omics approaches may present additional contributions to improve recombinant MP
798 production processes in *P. pastoris*.

799 Aiming to overcome the cellular burden caused by MP overexpression, researchers have
800 been driving their efforts towards the isolation and/or engineering of host cells, which have
801 proven to be efficient in many cases. In addition, codon usage optimizations have been shown to
802 be an effective strategy towards the improvement of MP expression but researchers should be
803 aware that synonymous mutations can affect protein function. The application of fusion partners
804 is helpful to increase MP solubility or to easily detect their expression levels and the advent of
805 transcriptional fusions show that particularly for solubility-enhancing tags, it seems that a
806 physical linkage between target MP and fusion may not be necessary for the desired effect, thus
807 simplifying the overall process.

808 Overall, the increasing understanding of MP biogenesis and the host physiological
809 response to MP recombinant production has allowed important advances in this field. However,
810 while it remains difficult to set general rules for a successful MP production process, the
811 information gathered in this review can help researchers with their own MP targets.

812

813 **Compliance with ethical standards:**

814 **Funding:** The authors acknowledge to the CICS-UBI projects Pest-
815 OE/SAU/UI0709/2014, UID/Multi/00709/2013 and the program COMPETE, Pest-
816 C/SAU/UI709/2011, financed by national funds through the FCT/MEC and when appropriate
817 co-financed by FEDER. CICS-UBI was also supported by FEDER funds through the POCI –
818 COMPETE 2020 – Operational Programme Competitiveness and Internationalisation in Axis I
819 – Strengthening research, technological development and innovation (Project POCI-01-0145-
820 FEDER-007491). This work was also developed within the scope of the project CICECO-
821 Aveiro Institute of Materials, FCT Ref. UID/CTM/50011/2019, financed by national funds
822 through the FCT/MCTES. The authors also acknowledge FCT for funding (Projects REFs:
823 EXPL/BBB478/BQB/0960/2012 and POCI-01-0145-FEDER-030840). Augusto Q. Pedro
824 acknowledges a doctoral fellowship (SFRH/BD/81222/2011) from FCT.

825

826 **Conflict of interest:** The authors declare that they have no conflict of interest.

827

828 **Ethical statement:** This article does not contain any studies with human participants or
829 animals performed by any of the authors.

830

831 **References:**

- 832 Abad S, Kitz K, Schreiner U, Hoermann A, Hartner F, Glieder A (2010) Real-time PCR-based
833 determination of gene copy numbers in *Pichia pastoris*. *Biotechnol J* 5 (4): 413-20.
- 834 Almo SC, Love JD (2014) Better and faster: improvements and optimization for mammalian
835 recombinant protein production. *Curr Opin Struct Biol* 26: 39-43.
- 836 André N, Cherouati N, Prual C, Steffan T, Zeder-Lutz G, Magnin T, Pattus F, Michel H,
837 Wagner R, Reinhart C (2006) Enhancing functional production of G protein-coupled receptors
838 in *Pichia pastoris* to levels required for structural studies via a single expression screen. *Protein*
839 *Sci* 15: 1115-26.
- 840 Andréll J, Tate CG (2013) Overexpression of membrane proteins in mammalian cells for
841 structural studies. *Mol Membr Biol* 30 (1): 52-63.
- 842 Angov E, Hillier CJ, Kincaid RL, Lyon JA (2008) Heterologous protein expression is enhanced
843 by harmonizing the codon usage frequencies of the target gene with those of the expression
844 host. *PLoS One* 3 (5): 2189-99.
- 845 Angov E (2011) Codon usage: nature's roadmap to expression and folding of proteins.
846 *Biotechnol J* 6 (6): 650-659.
- 847 Athey J, Alexaki A, Osipova E, Rostovtsev A, Santana-Quintero LV, Katneni U, Simonyan V,
848 Kimchi-Sarfaty C (2017) A new and updated resource for codon usage tables. *BMC*
849 *Bioinformatics* 18: 391.
- 850 Aw R, Polizzi KM (2013) Can too many copies spoil the broth? *Microb Cell Fact* 12: 128-37.
- 851 Backlund E, Ignatushchenko M, Larsson G (2011) Suppressing glucose uptake and acetic acid
852 production increases membrane protein overexpression in *Escherichia coli*. *Microb Cell Fact*
853 10: 35.
- 854 Bai J, Swartz DJ, Protasevich II, Brouillette CG, Harrell PM, Hildebrandt E, Gasser B,
855 Mattanovich D, Ward A, Chang G, Urbatsch IL (2011) A gene optimization strategy that
856 enhances production of fully functional P-glycoprotein in *Pichia pastoris*. *PLoS One* 6 (8):
857 22577-92.
- 858 Bandaranayake AD, Almo SC (2014) Recent advances in mammalian protein production. *FEBS*
859 *Lett* 588 (2): 253-260.
- 860 Baneyx F (1999) Recombinant protein expression in *Escherichia coli*. *Curr Opin Biotechnol* 10
861 (5): 411-421.
- 862 Baumgarten T, Schlegel S, Wagner S, Low M, Eriksson J, Bonde I, Herrgard MJ, Heipieper HJ,
863 Norholm MH, Slotboom DJ, de Gier JW (2017) Isolation and characterization of the *E. coli*
864 membrane protein production strain Mutant56(DE3). *Sci Rep* 7: 45089.
- 865 Baumruck AC, Tietze D, Steinacker LK, Tietze AA (2018) Chemical synthesis of membrane
866 proteins: a model study on the influenza virus B proton channel. *Chem Sci* 9: 2365-2375.

867 Bernaudat F, Frelet-Barrand A, Pochon N, Dementin S, Hivin P, Boutigny S, Rioux JB, Salvi D,
868 Seigneurin-Berny D, Richaud P, Joyard J, Pignol D, Sabaty M, Desnos T, Pebay-Peyroula E,
869 Darrouzet E, Vernet T, Rolland N (2011) Heterologous expression of membrane proteins:
870 choosing the appropriate host. *PLoS One* 6 (12): 29191-208.

871 Bollin F, Dechavanne V, Chevalet L (2011) Design of experiment in CHO and HEK transient
872 transfection condition optimization. *Protein Expr Purif* 78 (1): 61-68.

873 Brooks CL, Morrison M, Joanne Lemieux M (2013) Rapid expression screening of eukaryotic
874 membrane proteins in *Pichia pastoris*. *Protein Sci* 22 (4): 425-433.

875 Browne SM, Al-Rubeai, M (2007) Selection methods for high-producing mammalian cell lines.
876 *Trends Biotechnol* 25 (9): 425-32.

877 Byrne B (2015) *Pichia pastoris* as an expression host for membrane protein structural biology.
878 *Curr Opin Struct Biol* 32: 9-17.

879 Chaudhary S, Pak JE, Gruswitz F, Sharma V, Stroud RM (2012) Overexpressing human
880 membrane proteins in stably transfected and clonal human embryonic kidney 293S cells. *Nat*
881 *Protoc* 7 (3): 453-66.

882 Claassens NJ, Siliakus MF, Spaans SK, Creutzburg SCA, Nijse B, Schaap PJ, Quax TEF, van
883 der Oost J (2017) Improving heterologous membrane protein production in *Escherichia coli* by
884 combining transcriptional tuning and codon usage algorithms. *PLoS One* 12(9): e0184355.

885 Costa S, Almeida A, Castro A, Domingues L (2014) Fusion tags for protein solubility,
886 purification, and immunogenicity in *Escherichia coli*: the novel Fh8 system. *Front Microbiol* 5:
887 63.

888 Cregg JM, Cereghino JL, Shi J, Higgins DR (2000) Recombinant protein expression in *Pichia*
889 *pastoris*. *Mol Biotechnol* 16: 23-52.

890 Dilworth MV, Piel MS, Bettaney KE, Ma P, Luo J, Sharples D, Poyner DR, Gross SR, Moncoq
891 K, Henderson PJF, Miroux B, Bill RM (2018) Microbial expression systems for membrane
892 proteins. *Methods* 147: 3-39.

893 Drew DE, von Heijne G, Nordlund P, de Gier JW (2001) Green Fluorescent protein as an
894 indicator to monitor membrane protein overexpression in *Escherichia coli*. *FEBS Lett* 507 (2):
895 220-4.

896 Drew DE, Lerch M, Kunji E, Slotboom DJ, de Gier JW (2006) Optimization of membrane
897 protein overexpression and purification using GFP fusions. *Nat Methods* 3 (4): 303-313.

898 Elegheert J, Behiels E, Bishop B, Scott S, Woolley RE, Griffiths SC, Byrne EFX, Chang VT,
899 Stuart DI, Jones EY, Siebold C, Aricescu AR (2018) Lentiviral transduction of mammalian cells
900 for fast, scalable and high-level production of soluble and membrane proteins. *Nat Protoc* 13:
901 2991-3017.

902 Emmerstorfer-Augustin A, Wriessnegger T, Hirz M, Zellnig G, Pichler H (2019) Membrane
903 protein production in yeast: modification of yeast membranes for human membrane protein
904 production. In *Recombinant Protein Production in Yeast, Methods in Molecular Biology*,
905 Gasser B, Mattanovich D (Eds), Springer Nature, 1923: 265-285.

906 Fernández FJ, Vega MC (2016) Choose a suitable expression host: a survey of available protein
907 production platforms. *Advanced Technologies for protein complex production and*
908 *characterization*, Adv Exp Med Biol 896: 15-24.

909 Fuh MM, Steffen P, Schluter H (2016) Tools for the analysis and characterization of therapeutic
910 protein species. *Biosimilars* 6: 17-24.

911 Gan J, Ben-Nissan G, Arkind G, Tarnavsky M, Trudeau D, Garcia LN, Tawfik DS, Sharon M
912 (2017) Native mass spectrometry of recombinant proteins from crude cell lysates. *Ana Chem* 89
913 (8): 4398-4404.

914 Giacalone MJ, Gentile AM, Lovitt BT, Berkley NL, Gunderson CW, Surber MW (2006) Toxic
915 protein expression in *Escherichia coli* using a rhamnose-based tightly regulated and tunable
916 promoter system. *Biotechniques* 40: 355-64.

917 Gialama D, Kostelidou K, Michou M, Delivoria DC, Kolisis FN, Skretas G (2017)
918 Development of *Escherichia coli* strains that withstand membrane protein-induced toxicity and
919 achieve high-level recombinant membrane protein production. *ACS Synth Biol* 6 (2): 284-300.

920 Goehring A, Lee CH, Wang KH, Michel JC, Claxton DP, Bacongus I, Althoff T, Fischer S,
921 Garcia KC, Gouaux E (2014) Screening and large-scale expression of membrane proteins in
922 mammalian cells for structural studies. *Nat Protoc* 9 (11): 2574-85.

923 Gonçalves AM, Pedro AQ, Maia C, Sousa F, Queiroz JA, Passarinha LA (2013) *Pichia*
924 *pastoris*: A recombinant microfactory for antibodies and human membrane proteins. *J*
925 *Microbiol Biotechnol* 23 (5): 587-601.

926 Gould N, Hendy O, Papamichail D (2014) Computational tools and algorithms for designing
927 customized synthetic genes. *Front Bioeng Biotechnol* 2: 41.

928 Guerfal M, Ryckaert S, Jacobs PP, Jacobs PP, Ameloot P, Van Craenenbroeck K, Derycke R,
929 Callewaert N (2010) The HAC1 gene from *Pichia pastoris*: characterization and effect of its
930 overexpression on the production of secreted, surface displayed and membrane proteins. *Microb*
931 *Cell Fact* 9 (49): 2859-71.

932 Gul N, Linares DM, Ho FY, Poolman B (2014) Evolved *Escherichia coli* strains for amplified,
933 functional expression of membrane proteins. *J Mol Biol* 426 (1): 136-49.

934 Gustafsson C, Govindarajan S, Minshull J (2004) Codon bias and heterologous protein
935 expression. *Trends Biotechnol* 22 (7): 346-353.

936 Hacker DL, Balasubramanian S (2016) Recombinant protein production from stable mammalian
937 cell lines and pools. *Curr Opin Struct Biol* 38: 129-136.

938 Hardy D, Desuzinges Mandon E, Rothnie AJ, Jawhari A (2018) The yin and yang of
939 solubilization and stabilization for wild-type and full-length membrane protein. *Methods* 147:
940 118-125.

941 Hartmann L, Metzger E, Ottelard N, Wagner R (2017) Direct extraction and purification of
942 recombinant membrane proteins from *Pichia pastoris* protoplasts. *Methods Mol Biol* 1635: 45-
943 56.

944 Hattab G, Warschawski DE, Moncoq K, Miroux B (2015) *Escherichia coli* as host for
945 membrane protein structure determination: a global analysis. *Sci Rep* 5: 12097.

946 He M, He Y, Luo Q Wang M (2011) From DNA to protein: No living cells required. *Process*
947 *Biochem* 46: 615-20.

948 He Y, Wang K, Yan N (2014) The recombinant expression systems for structure determination
949 of eukaryotic membrane proteins. *Protein Cell* 5 (9): 658-72.

950 Henrich E, Hein C, Dotsch V, Bernhard F (2015) Membrane protein production in *Escherichia*
951 *coli* cell-free lysates. *Febs Lett* 589: 1713-22.

952 Hirz M, Richter G, Leitner E, Wriessnegger T, Pichler H (2013) A novel cholesterol-producing
953 *Pichia pastoris* strain is an ideal host for functional expression of human Na,K-ATPase $\alpha 3\beta 1$
954 isoform. *Appl Microbiol Biotechnol* 97: 9465-78.

955 Hsu M, Yu T, Chou C, Fu HY, Yang CS, Wang AH (2013) Using *Haloarcula marismortui*
956 bacteriorhodopsin as a fusion tag for enhancing and visible expression of integral membrane
957 proteins in *Escherichia coli*. *PLoS One* 8 (2): e56363.

958 Jacobs PP, Geysens S, Vervecken W, Contreras R, Callewaert N (2009) Engineering complex-
959 type N-glycosylation in *Pichia pastoris* using GlycoSwitch technology. *Nat Protoc* 4 (1): 58-70.

960 Jensen HM, Eng T, Chubukov V, Herbert RA, Mukhopadhyay A (2017) Improving membrane
961 protein expression and function using genomic edits. *Sci Rep* 7: 13030.

962 Kim YH, Youn H, Na J, Hong KJ, Kang KW, Lee DS, Chung JK (2015) Codon-optimized
963 human sodium iodide symporter (opt-hNIS) as a sensitive reporter and efficient therapeutic
964 gene. *Theranostics* 5 (1): 86-96.

965 Kopanic JL, Al-Mugotir M, Zach S, Das S, Grosely R, Sorgen PL (2013) An *Escherichia coli*
966 strain for expression of the connexin45 carboxyl terminus attached to the 4th transmembrane
967 domain. *Front Pharmacol* 4: 106.

968 Lantez V, Nikolaidis I, Rechenmann M, Vernet T, Noirclerc-Savoie M (2015) Rapid automated
969 detergent screening for the solubilization and purification of membrane proteins and complexes.
970 *Eng Life Sci* 15: 39-50.

971 Laukens B, De Wachter C, Callewaert N (2015) Engineering the *Pichia pastoris* N-
972 Glycosylation pathway using the GlycoSwitch technology. *Methods Mol Biol* 1321: 103-22.

973 Lee C, Kim J, Shin SG, Hwang S (2006) Absolute and relative quantification of plasmid copy
974 number in *Escherichia coli*. J Biotechnol 123: 273-80.

975 Lee JY, Chen H, Liu A, Alba BM, Lim AC (2017) Auto-induction of *Pichia pastoris* AOX1
976 promoter for membrane protein expression. Protein Expr Purif 137: 7-12.

977 Liu J, Srinivasan P, Pham DN, Rozovsky S (2012) Expression and purification of the membrane
978 enzyme selenoprotein K. Protein Expr Purif 86 (1): 27-34.

979 Löw C, Jegerschöld C, Kovermann M, Moberg M, Nordlund P (2012) Optimisation of over-
980 expression in *E. coli* and biophysical characterization of human membrane protein synaptogyrin
981 1. PLoS One 7 (6): 38244-57.

982 Luo J, Choulet J, Samuelson JC (2009) Rational design of a fusion partner for membrane
983 protein expression in *E. coli*. Protein Sci 18: 1735-44.

984 Lyons JA, Shahsavari A, Paulsen PA, Pedersen BP, Nissen P (2016) Expression strategies for
985 structural studies of eukaryotic membrane proteins. Curr Opin Struct Biol 38: 137-144.

986 Ma C, Hao Z, Huysmans G, Lesiuk A, Bullough P, Wang Y, Bartlam M, Phillips SE, Young
987 JD, Goldman A, Baldwin SA, Postis VL (2015) A versatile strategy for production of
988 membrane proteins with diverse topologies: application to investigation of bacterial homologues
989 of human divalent metal ion and nucleoside transporters. PLoS One 10 (11): e10143010.

990 Mancina F, Patel SD, Rajala MW, Scherer PE, Nemes A, Schieren I, Hendrickson WA, Shapiro
991 L (2004) Optimization of protein production in mammalian cells with a coexpressed fluorescent
992 marker. Structure 12: 1355-60.

993 Marino J, Hohl M, Seeger MA, Zerbe O, Geertsma ER (2015) Bicistronic mRNA to enhance
994 membrane protein overexpression. J Mol Biol 427 (4): 943-54.

995 Marino J, Holzhter K, Kuhn B, Geertsma ER (2017) Efficient screening and optimization of
996 membrane protein production in *Escherichia coli*. Methods Enzymol 594: 139-164.

997 Martins LM, Pedro AQ, Oppolzer D, Sousa F, Queiroz JA, Passarinha LA (2015) Enhanced
998 biosynthesis of plasmid DNA from *Escherichia coli* VH33 using Box-Behnken design
999 associated to aromatic amino acids pathway. Biochem Eng J 98: 117-26.

1000 Massey-Gendel E, Zhao A, Boulting G, Kim H-Y, Balamotis MA, Nakamoto RK, Bowie JU
1001 (2009) Genetic selection system for improving recombinant membrane protein expression in *E.*
1002 *coli*. Protein Sci 18: 372-83.

1003 Mauro VP (2018) Codon optimization in the production of recombinant biotherapeutics:
1004 potential risks and considerations. BioDrugs 32 (1): 69-81.

1005 McKenzie EA, Abbott WM (2018) Expression of recombinant proteins in insect and
1006 mammalian cells. Methods 147: 40-49.

1007 Midgett CR, Madden DR (2007) Breaking the bottleneck: Eukaryotic membrane protein
1008 expression for high-resolution structural studies. J Struct Biol 160: 265-74.

1009 Miroux B, Walker JE (1996) Over-production of proteins in *Escherichia coli*: mutant hosts that
1010 allow synthesis of some membrane proteins and globular proteins at high-levels. *J Mol Biol* 260
1011 (3): 289-298.

1012 Mizrachi D, Chen Y, Liu J, Peng HM, Ke A, Pollack L, Turner RJ, Auchus RJ, DeLisa MP
1013 (2015) Making water-soluble integral membrane proteins in vivo using an amphipatic protein
1014 fusion strategy. *Nat Commun* 6: 6826.

1015 Murata Y, Watanabe T, Sato M, Momose Y, Nakahara T, Oka S, Iwahashi H (2003) Dimethyl
1016 sulfoxide exposure facilitates phospholipid biosynthesis and cellular membrane proliferation in
1017 yeast cells. *J Biol Chem* 278: 33185-33193.

1018 Nannenga BL, Baneyx F (2011) Reprogramming chaperone pathways to improve membrane
1019 protein expression in *Escherichia coli*. *Protein Sci* 20: 1411-20.

1020 Narayanan A, Ridilla M, Yernool DA (2011) Restrained expression, a method to overproduce
1021 toxic membrane proteins by exploiting operator-repressor interactions. *Protein Sci* 20 (1): 51-
1022 61.

1023 Nieuwkoop T, Claassens NJ, van der Oost J (2019) Improved protein production and codon
1024 optimization analyses in *Escherichia coli* by bicistronic design. *Microb Biotechnol* 12 (1): 173-
1025 179.

1026 Nji E, Chatzikyriakidou Y, Landreh M, Drew D (2018) An engineered thermal-shift screen
1027 reveals specific lipid preferences of eukaryotic and prokaryotic membrane proteins. *Nat*
1028 *Commun* 9: 4253.

1029 Nordén K, Agemark M, Danielson JA, Alexandersson E, Kjellbom P, Johanson U (2011)
1030 Increasing gene dosage greatly enhances recombinant expression of aquaporins in *Pichia*
1031 *pastoris*. *BMC Biotechnol* 11: 47-59.

1032 Nørholm MH, Toddo S, Virkki MT, Light S, von Heijne G, Daley DO (2013) Improved
1033 production of membrane proteins in *Escherichia coli* by selective codon substitutions. *FEBS*
1034 *Lett* 587 (15): 2352-2358.

1035 Oberg F, Ekvall M, Nyblom M, Backmark A, Neutze R, Hedfalk K (2009) Insight into factors
1036 directing high production of eukaryotic membrane proteins; production of 13 human AQPs in
1037 *Pichia pastoris*. *Mol Membr Biol* 26 (4): 215-27.

1038 Ohsfeldt E, Huang S, Baycin-Hizal D, Kristoffersen L, Le TM, Li E, Hristova K, Betenbaugh
1039 MJ (2012) Increased expression of the integral membrane proteins EGFR and FGFR3 in anti-
1040 apoptotic Chinese Hamster Ovary cell lines. *Biotechnol Appl Biochem* 59: 155-62.

1041 Oliveira C, Domingues L (2018) Guidelines to reach high-quality purified recombinant proteins.
1042 *Appl Microbiol Biotechnol* 102 (1): 81-92.

1043 Pandey A, Shin K, Patterson RE, Liu X, Rainey JK (2016) Current strategies for protein
1044 production and purification enabling membrane protein structural biology. *Biochem Cell Biol*
1045 94, 507-527.

1046 Parret AH, Besir H, Meijers R (2016) Critical reflections on synthetic gene design for
1047 recombinant protein expression. *Curr Opin Struct Biol* 38: 155-162.

1048 Pedro AQ, Martins LM, Dias JM, Bonifácio MJ, Queiroz JA, Passarinha LA (2015) An
1049 artificial neural network for membrane-bound catechol-*O*-methyltransferase biosynthesis with
1050 *Pichia pastoris* methanol-induced cultures. *Microb Cell Fact* 14: 113-27.

1051 Popot JL (2018) Membrane proteins in Aqueous solution, From detergents to amphipols.
1052 Springer International Publishing.

1053 Puigbò P, Guzmán E, Romeu A, Garcia-Vallvé S (2007) OPTIMIZER: a web server for
1054 optimizing the codon usage of DNA sequences. *Nucleic Acid Res* 35 (2): 126-31.

1055 Puigbo P, Bravo IG, Garcia-Vallve S (2008) CAIcal: a combined set of tools to assess codon
1056 usage adaptation. *Biol Direct* 3: 38.

1057 Quax TE, Claassens NJ, Soll D, van der Oost J (2015) Codon bias as a means to fine-tune gene
1058 expression. *Mol Cell* 59 (2): 149-161.

1059 Rahman M, Ismat F, McPherson MJ, Baldwin SA (2007) Topology-informed strategies for the
1060 overexpression and purification of membrane proteins. *Mol Membr Biol* 24: 407-418.

1061 Rajesh S, Knowles T, Overduin M (2011) Production of membrane proteins without cells or
1062 detergents. *N Biotechnol* 28 (3): 250-54.

1063 Ramón A, Marín M (2011) Advances in the production of membrane proteins in *Pichia*
1064 *pastoris*. *Biotechnol J* 6: 700-6.

1065 Raynal B, Lenormand P, Baron B, Hoos S, England P (2014) Quality assessment and
1066 optimization of purified protein samples: why and how? *Microb Cell Fact* 13: 180.

1067 Rosano GL, Ceccarelli EA (2014) Recombinant expression in *Escherichia coli*: advances and
1068 challenges. *Front Microbiol* 5: 172.

1069 Saladi SM, Javed N, Muller A, Clemons WM Jr (2018) A statistical model for improved
1070 membrane protein expression using sequence-derived features. *J Biol Chem* 293 (13): 4913-
1071 4927.

1072 Schlegel S, Lofblom J, Lee C, Hjelm A, Klepsch M, Strous M, Drew D, Slotboom DJ, de Gier
1073 JW (2012) Optimizing membrane protein overexpression in the *Escherichia coli* Lemo21
1074 (DE3). *J Mol Biol* 423: 648-59.

1075 Schlegel S, Genevaux P, de Gier J (2017) Isolating *Escherichia coli* strains for recombinant
1076 protein production. *Cell Mol Life Sci* 74 (5): 891-908.

1077 Sharp PM, Li WH (1987) The codon adaptation index – a measure of directional synonymous
1078 codon usage bias, and its potential applications. *Nucleic Acids Res* 15 (3): 1281-1295.

1079 Shiroishi M, Kobayashi T, Ogasawara S, Tsujimoto H, Ikeda-Suno C, Iwata S, Shimamura T
1080 (2011) Production of the stable human histamine H₁ receptor in *Pichia pastoris* for structural
1081 determination. *Methods* 55 (4): 281-286.

1082 Shukla S, Schwartz C, Kapoor K, Kouanda A, Ambudkar SV (2012) Use of baculovirus
1083 BacMam vectors for expression of ABC drug transporters in mammalian cells. *Drug Metab*
1084 *Dispos* 40 (2): 304-12.

1085 Skretas G, Makino T, Varadaraiyan N, Pogson M, Georgiou G (2012) Multi-copy genes that
1086 enhance the yield of mammalian G protein-coupled receptors in *Escherichia coli*. *Metab Eng* 14
1087 (5): 591-602.

1088 Snijder HJ, Hakulinen J (2016) Membrane protein production in *E. coli* for applications in drug
1089 discovery. *Adv Exp Med Biol* 896: 59-77.

1090 Talmont F, Sidobre S, Demange P, Milon A, Emorine LJ (1996) Expression and
1091 pharmacological characterization of the human mu-opioid receptor in the methylotrophic yeast
1092 *Pichia pastoris*. *FEBS Lett* 394: 268-272.

1093 Van der Rest ME, Kamminga AH, Nakano A, Anraku Y, Poolman B, Konings WN (1995) The
1094 plasma membrane of *Saccharomyces cerevisiae*: structure, function and biogenesis. *Microbiol*
1095 *Rev* 59: 304-322.

1096 Vogl T, Thallinger GG, Zellnig G, Drew D, Cregg JM, Glieder A, Freigassner M (2014)
1097 Towards improved membrane protein production in *Pichia pastoris*: General and specific
1098 transcriptional response to membrane protein overexpression. *N Biotechnol* 31 (6): 538-552.

1099 Wagner S, Bader ML, Drew D, de Gier J (2006) Rationalizing membrane protein
1100 overexpression. *Trends Biotechnol* 24 (8): 364-71.

1101 Wagner S, Klepsch MM, Schlegel S, Appel A, Draheim R, Tarry M, Hogbom M, van Wijk KJ,
1102 Slotboom DJ, Persson JO, de Gier JW (2008) Tuning *Escherichia coli* for membrane protein
1103 overexpression. *PNAS* 105 (38): 14371-14376.

1104 Welch M, Villalobos A, Gustafsson C, Minshull J (2011) Designing genes for successful
1105 protein expression. *Methods Enzymol* 498: 43-66.

1106 Wen Z, Boddicker MA, Kaufhold RM, Khandelwal P, Durr E, Qiu P, Lucas BJ, Nahas DD,
1107 Cook JC, Touch S, Skinner JM, Espeseth AS, Przysiecki CT, Zhang L (2016) Recombinant
1108 expression of *Chlamydia trachomatis* major outer membrane protein in *E. coli* outer membrane
1109 as a substrate for vaccine research. *BMC Microbiol* 16: 165.

1110 Zhang G, Annan RS, Carr SA, Neubert TA (2010) Overview of peptide and protein analysis by
1111 mass spectrometry. *Curr Protoc Protein Sci* 62 (1): 16.1.1-16.1.30.

1112 Zheng X, Dong S, Zheng J, Li D, Li F, Luo Z (2014) Expression, stabilization and purification
1113 of membrane proteins via diverse protein synthesis systems and detergents involving cell-free
1114 associated with self-assembly peptide surfactants. *Biotech Adv* 32: 564-74.

1115 Zuo X, Li S, Hall J, Mattern MR, Tran H, Shoo J, Tan R, Weiss SR, Butt TR (2005) Enhanced
1116 expression and purification of membrane proteins by SUMO fusion in *Escherichia coli*. J Struct
1117 Funct Genomics 6: 103-11.

Tables:

Table 1 - Major advantages, limitations and general characteristics of recombinant membrane protein expression systems.

<i>Host</i>	<i>Advantages</i>	<i>Drawbacks</i>	<i>Other characteristics</i>	<i>References</i>
<i>Escherichia coli</i> Gram-negative bacterium	Inexpensive; Rapid generation of expression plasmids; Fast growth; Easy scale up; Simple culture requirements.	Endotoxin; Inclusion body formation; Inefficient protein secretion; Many MP do not fold properly; Lack of efficient PTM; Unable to efficiently express proteins larger than 120 kDa.	Specific strains (e.g. Lemo21) or introduction of solubility tags may improve MP expression. Inner membrane and the inner leaflet of the outer membrane are mainly composed by phosphatidylethanolamine followed by phosphatidylglycerol and few cardiolipin and the outer leaflet of the outer membrane is highly enriched in Lipopolysaccharide.	(Bernaudat et al 2011; Midgett et al 2007; McMorran et al 2014; Fernández and Vega, 2016)
<i>Pichia pastoris</i> Methylotrophic yeast; GRAS organism	Efficient protein secretion with low levels of endogenous proteins. Capable of performing many PTM; Low cost of culture media; Industry-scale fermentation.	Glycosylation pattern different from mammalian; Intracellular recovery of large amount of cells may require specific equipment (French-press); High oxygen demand.	Improved glyco-engineered strains obtained using the GlycoSwitch® technology; Wide range of genetic tools, plasmids, strains and promoters available; The preference for the respiratory growth allow to be cultivated at high cell densities. Plasma membrane composed of phospholipids, sterols (ergosterol) and sphingolipids (inositol).	(Gonçalves et al 2013; Laukens et al 2015; Marredy et al 2011; Pedro et al 2015)
Insect cells Baculovirus-infected cells	More native environment than yeast; More compatible with eukaryotic MP because of similar codon usage rules than <i>E. coli</i> or <i>P. pastoris</i> ; Well-established protocols; Good secretion.	Cost; Non-native glycosylation and lipid environment; Cell lysis; Some of the PTM are not identical to those found in mammalian; Long production time; Relative high media costs.	Used for MP expression as a compromise between bacterial and mammalian systems. Viral infection promote cell lysis and may lead to proteolysis of target protein.	(Bernaudat et al 2013; Midgett et al 2007)
Mammalian cells Stable integration and transient transfection	Proper folding; Stable/transient folding; Native lipid environment and post-translational pathways.	High media costs; Slow growth rates; Low expression; Viral infection; Cost; Higher technical requirements.	For particular targets, may be the only expression system able to express a given MP in a functional and properly folded state. Cholesterol present in membranes may be essential for the functionality of certain MP.	(Midgett et al 2007; Andréll et al 2013)
Cell-Free expression	Short time reaction; Manipulation of reaction conditions allow to control conveniently the PTM; Plasmid or DNA can be directly used for protein expression; Special proteins can be expressed with a composition of non-natural amino acids.	High costs, Low protein production rates; Insufficiency of PTM is a bottleneck to obtain complex proteins in a functional form.	May be based in prokaryotic or eukaryotic CF systems; MP may be produced co-translationally in artificial membrane environments.	(Rajesh et al 2011; Proverbio et al 2013; Zheng et al 2014)

Table 2 – Critical assessment of major parameters affecting the upstream stage of recombinant MP structural biology projects for a good decision-making process.

Parameter	<i>Escherichia coli</i>	<i>Pichia pastoris</i>	Mammalian cell lines		Baculovirus-infected Insect cells
			Transiently transfected	Stable clones	
Gene Dosage	Preference: Plasmid-based system with medium – high PCN	Mixed results, screening is advisable; higher gene dosage can increase yield	Usually favored by high gene dosage		MOI affects expression yields
Codon optimization	Advisable testing for heterologous targets; “harmonized” codons often leads to outstanding improvements				
Cost	Very low	Low	Very high		High
Ease of manipulation/Labor intensive	High/Low	High/Low	Low/High	Low/Very High	Low/High
Scalability	Very Good	Very Good	Moderate		Moderate
Timescale	Days	Week	Days/Week	Lengthy (Months)	Weeks
Membrane Protein features	Glycosylation	Low/Absent	High ¹		High
	Other PTM	Bad	Good		Good
	Lipid composition	Bad	Good ²		Good
	Specific organelles requirements (e.g. mitochondria)	Bad	Bad		Good
	Molecular Weight	Limited	Good		Good
	Protein productivity	Good	Very Good		Bad
Observations	The source from which more MP structures were solved ³	Viable alternative to mammalian and insect cells for obtaining low cost and high yield of MP	The most complete for human MP expression, greatly exemplified by SERT ⁴		Applied as a compromise between bacteria and mammalian cell lines

Legend: – ¹Using the GlycoSwitch® technology (Laukens et al 2015); ²Humanized pathway (Hirz et al 2013); ³Pandey et al 2016; ⁴Andréll et al 2013.

Figures:

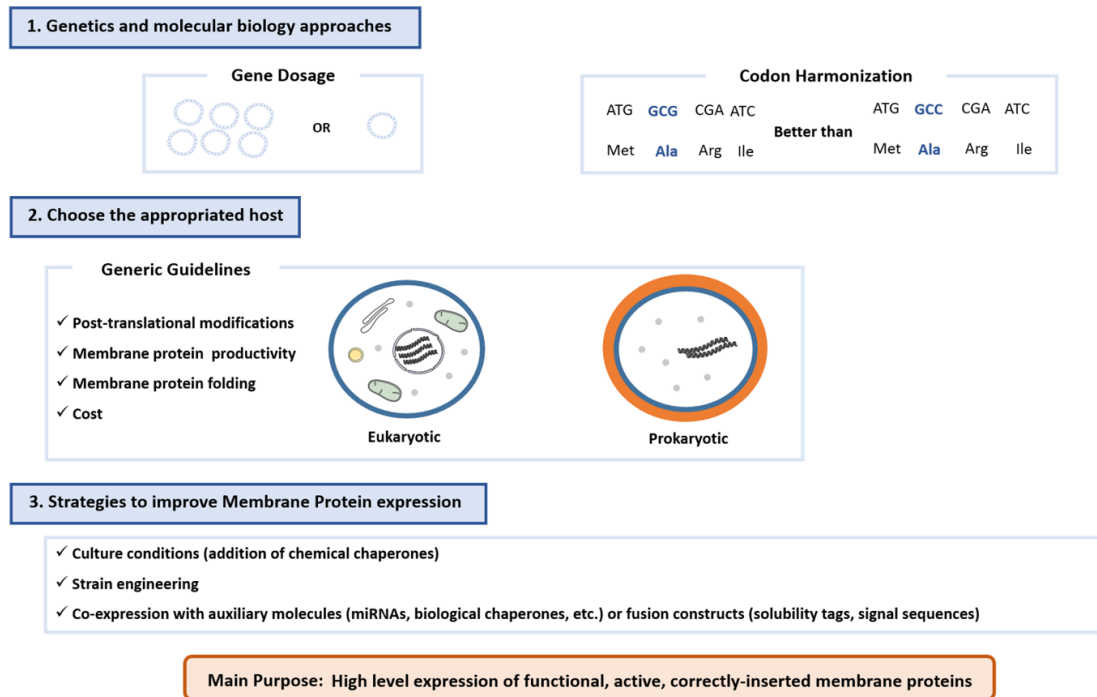


Figure 1 – Overview of the topics included in this review amenable to optimization and, thus, relevant for obtaining a successful strategy for recombinant MP biosynthesis.

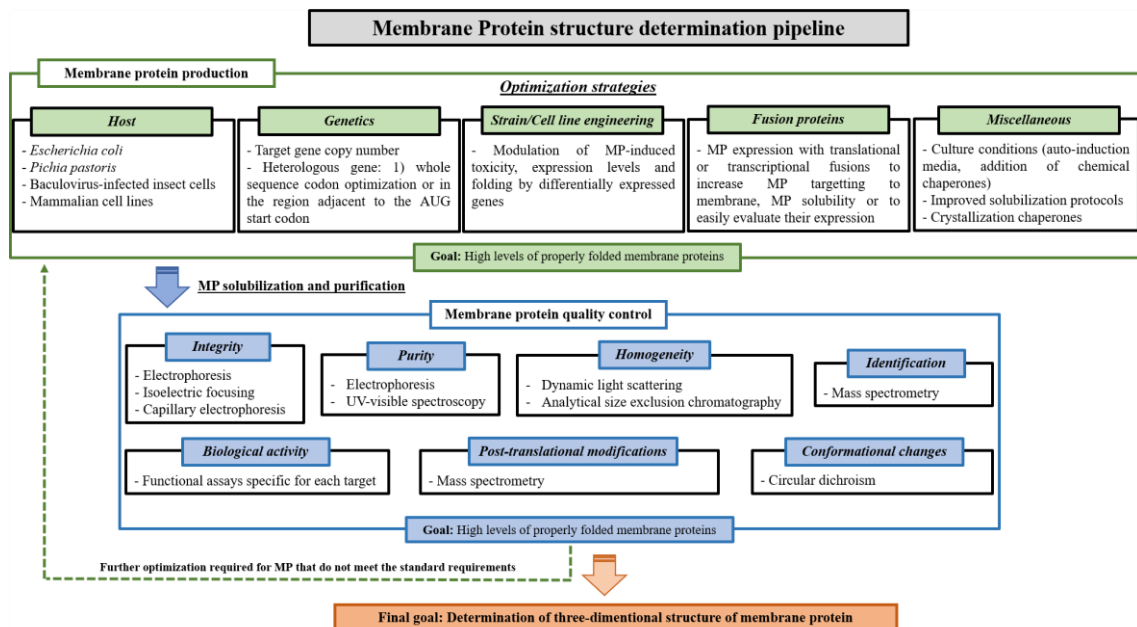


Figure 2 – Schematic diagram of MP structure determination pipeline focusing relevant parameters to optimize their upstream stage and techniques used to protein quality control.