# Understudied Proteins and Understudied Functions in the Model Bacterium Bacillus subtilis – a Major Challenge in Current Research

Jörg Stülke[1], Dennis Wicke[1], Janek Meißner[1], Robert Warneke[1], and Christoph Elfmann[1]

[1]Georg-August-Universitat Gottingen Abteilung fur Allgemeine Mikrobiologie

February 2, 2023

## Abstract

Model organisms such as the Gram-positive bacterium *Bacillus subtilis* have been studied intensively for decades. However, even for such model organisms no function has been identified for about one fourth of all proteins. It has recently been realized that such understudied proteins as well as poorly studied functions set a limitation to our understanding of the requirements for cellular life, and the Understudied Proteins Initiative has been launched. Obviously, poorly studied proteins that are strongly expressed, are likely to be most important to the cell and should therefore have priority in further studies. Since the functional analysis of unknown proteins can be extremely laborious, a minimal knowledge is required prior to targeted functional studies. In this review, we discuss strategies to obtain such a minimal annotation, e.g. from global interaction, expression or localization studies. We present a set of 41 highly expressed and poorly studied proteins of *B. subtilis*. Several of these proteins are thought or known to bind RNA and/or the ribosome, some may may control the metabolism of *B. subtilis*, and another subset of particularly small proteins may act as regulatory elements to control the expression of downstream genes. Moreover, we discuss the challenges of poorly studied functions with a focus on RNA-binding proteins, amino acid transport and the control of metabolic homeostasis. The identification of the functions of the selected proteins will strongly advance our knowledge on *B. subtilis*, but also on other organisms since many of the proteins are conserved in many groups of bacteria.

## Review

**Understudied Proteins and Understudied Functions in the Model Bacterium *Bacillus subtilis* – a Major Challenge in Current Research**

**Dennis Wicke,# Janek Meißner,# Robert Warneke, # Christoph Elfmann, and Jörg Stülke***

*Georg-August-University Göttingen, Department of General Microbiology, Grisebachstr. 8, D-37077 Göttingen, Germany*

# These authors contributed equally to this work.

* Corresponding author

Department of General Microbiology, Georg-August-University Göttingen, Grisebachstr. 8, 37077, Göttingen, Germany

Phone: +49-551-393781; Fax: +49-551-393808

*jstuelk@gwdg.de*

**Running title: Understudied proteins in *Bacillus subtilis***

**Keywords: RNA-binding proteins, protein-protein interaction, protein-RNA interaction, ribosome, metabolism, amino acid transport**

## ABSTRACT

Model organisms such as the Gram-positive bacterium *Bacillus subtilis* have been studied intensively for decades. However, even for such model organisms no function has been identified for about one fourth of all proteins. It has recently been realized that such understudied proteins as well as poorly studied functions set a limitation to our understanding of the requirements for cellular life, and the Understudied Proteins Initiative has been launched. Obviously, poorly studied proteins that are strongly expressed, are likely to be most important to the cell and should therefore have priority in further studies. Since the functional analysis of unknown proteins can be extremely laborious, a minimal knowledge is required prior to targeted functional studies. In this review, we discuss strategies to obtain such a minimal annotation, e.g. from global interaction, expression or localization studies. We present a set of 41 highly expressed and poorly studied proteins of *B. subtilis* . Several of these proteins are thought or known to bind RNA and/or the ribosome, some may may control the metabolism of *B. subtilis,* and another subset of particularly small proteins may act as regulatory elements to control the expression of downstream genes. Moreover, we discuss the challenges of poorly studied functions with a focus on RNA-binding proteins, amino acid transport and the control of metabolic homeostasis. The identification of the functions of the selected proteins will strongly advance our knowledge on *B. subtilis* , but also on other organisms since many of the proteins are conserved in many groups of bacteria.

## 1 INTRODUCTION

Several decades of biochemical research have resulted in the elucidation of functions of a large number of proteins. Yet, we are still far from a comprehensive understanding of the role of all proteins of any single organism. This is illustrated by the artificial minimal organism*Mycoplasma mycoides* JCVI-syn3A. With only 452 protein-coding genes, this organism currently defines the lower limit of the protein complement for an independently viable bacterium (Hutchison et al., 2016; Breuer et al., 2019). Yet, about one third of these proteins still have no known function, indicating major gaps in our knowledge even for the most simple cells (Hutchison et al., 2016; Pedreira et al., 2022 a). These gaps are caused by a focus on a limited set of intensively studied proteins for which more and more knowledge accumulates. On the other hand, the biological function of a significant number of proteins remains not at all or only poorly understood.

Recently, it has been proposed to close this gap of knowledge by launching the Understudied Proteins Initiative (Kustatscher et al., 2022 a, b). The functional analysis of unknown proteins can be highly laborious and poorly rewarding. This is exemplified by a European/ Japanese initiative to functionally identify all unknown proteins of the Gram-positive model bacterium *Bacillus subtilis* after the completion of the genome sequence. While tremendous human and financial resources went into this project, functions have not been identified for more than a handful of so far unknown proteins. This results from the repeated investigation of a defined set of phenotypic analyses under a defined set of conditions. Typically, these phenotypes and conditions have been studied intensively before, so that only little new knowledge can be expected. As a conclusion, at least a minimal amount of molecular/ functional annotation is required as the starting point to identify the function of so far unknown proteins. This minimal knowledge could cover expression over a wide range of conditions, the similarity of phenotypes to known phenotypes, or the association of unknown proteins with proteins of known function, with RNA or other biomolecules. Finally, the identification of gaps in our knowledge and the goal-driven investigation of these functions may help to unravel the functions of poorly studied proteins. Another prerequisite for closing the annotation gap is the integration of all available information in intuitively accessible databases.

We are interested in the model organism *B. subtilis* . As a model organism for differentiation and workhorse in biotechnology, *B. subtilis* is one of the most intensively studied organisms. However, the function of about 25% of the proteins (about 1,000 proteins) encoded by the *B. subtilis* genome is still unknown or only very poorly understood (Pedreira et al., 2022). In this review, we give an overview on the strategies used to

get initial minimal annotation for so far unknown proteins, we present a set of highly expressed unknown proteins that should be studied with highest priority, and we define and discuss fields of research that still have many open questions, *i.e* ., RNA-binding proteins, amino acid transport, and the control of metabolic homeostasis.

## 2 STRATEGIES TO OBTAIN MINIMAL ANNOTATION FOR UNKNOWN PROTEIN OF *B. SUBTILIS*

As outlined above, some information is required to start the identification of unknown proteins. In the past years, such starting information has been obtained for many unknown proteins. A major breakthrough in this respect was the global transcriptome analysis of this organism under 104 different conditions (Nicolas et al., 2012). Based on this analysis, 180 proteins could be newly categorized as sporulation proteins (Pedreira et al., 2022b).

Interaction studies might provide more insights into the function of unknown proteins – based on the principle "guilty by association". Indeed, a large scale, global interactome analysis with the minimal organism *Mycoplasma pneumoniae* provided important clues for the function of many proteins. As an example, the completely unknown protein MPN530 interacts with multiple subunits of the RNA polymerase, including the house-keeping sigma factor, suggesting that this protein plays a role in transcription (O'Reilly et al., 2020; Elfmann et al., 2022). A similar attempt has recently also been performed for *B. subtilis* , and many so far unknown proteins were found interacting with well-characterized proteins or protein complexes, such as the ribosome (O'Reilly et al., 2023). Similarly, global approaches to identify protein-RNA or protein-metabolite complexes have been established; however, they have not yet been applied to *B. subtilis* (Chihara et al., 2022; Link et al., 2013). In addition to global analyses, protein-specific interaction studies can also provide important functional information about previously unknown proteins. This is the case for the c-di-AMP binding DarB protein, which was shown to interact with and thereby to control the activity of the (p)ppGpp synthetase/ hydrolase Rel and the pyruvate carboxylase PycA in *B. subtilis*(Krüger et al., 2021a; Krüger et al., 2022).

Phenotypic profiling is another way to get at least some initial information for unknown proteins. Since the experimental study of mutant phenotypes under a variety of experimental conditions often fails to provide clues (see above), global analyses of mutant properties might be an important way to obtain initial information on unknown proteins. For example, transcriptome or proteome profiling of changes in a mutant and a comparison to established catalogues may yield functional insights. Such an approach has already been successfully used to identify the mode of action of unknown antibacterial compounds in *B. subtilis*(Senges et al., 2021; Senges et al., 2022). Proteomic profiling may prove to be especially useful for the functional analysis of unknown ribosome-interacting proteins, since extensive catalogued information of proteome changes under conditions of ribosome perturbation is available for *B. subtilis* (Senges et al., 2021).

While the strategies described above all depend on a gene- or protein-based investigation, the analysis of poorly studied functions may also provide novel functional information. While many areas of research have been exhaustively addressed in this model bacterium, there still remain some functions that are uncharted territories of research. Among these functions are RNA-binding proteins, the transport of amino acids, the control of metabolic homeostasis (see below), some activities in cell wall biosynthesis, or the detoxification of toxic metabolites (Reuss et al., 2016). The purposeful investigation of these fields aimed at identifying the responsible proteins indeed uncovers functional information for previously unknown proteins. This is the case for the undecaprenyl-phosphate transporter UptA, which was very recently identified in a well-designed transposon mutagenesis screen (Roney & Rudner, 2022). This was the last missing piece in the part list of cell envelope biogenesis in *B. subtilis* . Similarly, several complementary approaches aimed at the identification of a protein responsible for serine uptake pointed to the so far unknown protein YbeC. Subsequently, YbeC was also shown to be the major glutamate transporter of *B. subtilis.* Accordingly, the protein was renamed AimA for amino acid importer A (Klewing et al., 2020; Krüger et al., 2021b). The investigation of proteins involved in metabolite damage control revealed that the YqeK protein may act in the degradation of toxic side-products of the nicotinamide-nucleotide adenylyltransferase NadD by virtue of its

versatile diphosphatase activity (Haas et al., 2022). Other examples for toxic by-products of metabolism are 4-phosphoerythronate, a by-product of erythrose-4-phosphate oxidation in the pentose phosphate pathway that inhibits the phosphogluconate dehydrogenase GndA, and 5-oxoproline, an unavoidable damage product formed spontaneously from glutamine. These harmful metabolites are detoxified by the GTPase CpgA, which moonlights in dephosphorylation of 4-phosphoerythronate, and by the 5-oxoprolinase PxpABC, respectively (Sachla & Helmann, 2019; Niehaus et al., 2017). It has recently become obvious that the prevention of metabolite damage is very important for the viability of any living cell. The limited knowledge on these mechanisms is an important bottleneck in all genome reduction projects (Reuß et al., 2017).

A final prerequisite for the identification of unknown functions is the integration of all possible available information. Even small pieces of information that by themselves may not prove to be very useful can help to get a deeper understanding if brought into an appropriate context. This annotation information can cover expression data, interaction data, gene regulation, the control of protein activities, localization data and many other types of data (see Fig. 1). An example for the value of data integration is one of the very few remaining unknown essential proteins, YlaN. This gene is essential under standard conditions (complex medium) but becomes dispensable if iron is added to the medium, suggesting a role in the control of iron homeostasis (Peters et al., 2016). In addition, two independent studies identified a physical interaction of the protein with the key regulator of iron homeostasis, Fur (de Jong et al., 2021; O'Reilly et al., 2023). Bringing this information together immediately supports the idea that the interaction is meaningful and that YlaN might control the activity of the Fur regulator. This hypothesis may then be validated in specifically designed experiments. Such an integration of all available information on the genes and proteins is provided in the database *Subti* Wiki which is the major reference tool of the *B. subtilis* research community (Fig. 1; Pedreira et al., 2022). Other important online tools that help developing hypotheses are the protein interaction and association database STRING (Sklarczyk et al., 2023), the FlaGs webserver that allows to interrogate conserved gene organization which is often an indication of functional association between proteins (Saha et al. 2021), as well as the UniProt and COG databases that provide information on protein functions that may have been identified in other, otherwise potentially overlooked (Galperin et al., 2021; UniProt Consortium, 2023).

## 3 A SET OF HIGHLY EXPRESSED UNKNOWN PROTEINS IN *B. SUBTILIS*

A recent proteomic study with *B. subtilis* revealed that essential proteins are highly overrepresented in the proteome (Reuß et al., 2017). Although they account only for only 6% of all *B. subtilis* genes, the essential genes use 57% of the translation capacity. This reflects the importance of these proteins for the cell. On the other hand, unknown and poorly studied proteins, which correspond to 25% of the genes, only use 3% of the translation capacity under standard conditions. It is very likely that many of the unknown proteins are only needed under very special conditions. This special-purpose demand combined with the low expression also explains why the function of such proteins has never been identified. On the other hand, there is a set of 41 proteins of unknown function that are highly expressed under most conditions (Table 1). Among these are the two essential proteins YlaN and YneF. It is tempting to speculate that these highly expressed unknown proteins are important for the physiology of *B. subtilis* even under standard conditions. Interestingly, there are three pairs of paralogous proteins on the list (YqhY/ YloU, YabR/ YugI, YtxH/ YhaH). Moreover, several of the proteins have been detected in a recent global in vivo interaction study with *B. subtilis* (see Table 1). The complete set of highly expressed unknown B. subtilis proteins can be easily accessed in the database SubtiWiki (http://www.subtiwiki.uni-goettingen.de/v4/category?id=SW.6.7; Pedreira et al., 2022b; see Fig. 1B).

Of the 41 proteins, several have or may have RNA-binding activity, among them the most strongly expressed and highly conserved unknown protein YqeY. This protein contains a domain also present in the GatB subunit of the glutamyl-tRNA amidotransferase subunit suggesting that YqeY might also have tRNA amino acid amidase activity. The YtpR protein possesses an tRNA-binding domain that is also present in the beta subunit of the phenylalanine tRNA synthetase. YtpR physically interacts with GatB (O'Reilly et al., 2023) suggesting that it might facilitate the interaction between the Glu-loaded tRNA and the glutamyl-tRNA

4

amidotransferase GatABC. The paralogous YabR and YugI proteins possess an RNA-binding S1 domain. Both proteins interact with the small subunit of the ribosome (O'Reilly et al., 2023) indicating their involvement in translation. This may also be the case for YrzB, which physically interacts with multiple ribosomal proteins (O'Reilly et al., 2023). The YlxR and KhpA proteins were found to bind RNA in *Clostrioides difficile* (Lamm-Schmidt et al., 2021). Finally, the YlbN protein is conserved in most bacteria and plant chloroplasts. The orthologous chloroplast and *E. coli* proteins are required for the accumulation of 23S rRNA (Yang et al., 2016) even though the molecular activity of the protein remains unknown.

Several of the highly expressed unknown proteins are likely involved in the control of metabolism. The YjlC protein is encoded in an operon with the NADH dehydrogenase, and the two proteins interact physically (O'Reilly et al., 2023) (see Fig. 1A). It is tempting to speculate that YjlC somehow controls the activity of NAD dehydrogenase. Indeed, both proteins are required for genetic transformation (Koo et al., 2017) indicating that they perform a joint function. The YlaN protein interacts with the key regulator of iron homeostasis (O'Reilly et al., 2023), and the normally essential gene becomes dispensable at high iron concentrations (Peters et al., 2016). Thus, YlaN may control iron homeostasis via Fur. The paralogous YqhY and YloU proteins are encoded in operons with genes required for complementary aspects of fatty acid acquisition, either biosynthesis or fatty acid phosphorylation. The *yqhY* gene is quasi-essential and the cells respond to its inactivation with the accumulation of suppressor mutations in the subunits of acetyl-CoA carboxylase (Tödter et al., 2017). Thus, these two proteins may control different aspects of lipid biosynthesis. The strength of initial protein-protein interaction information is demonstrated by the YneR protein which interacts with the PdhA and PdhB subunits of pyruvate dehydrogenase. Targeted experimental studies revealed that YneR acts as an inhibitor of pyruvate dehydrogenase activity. The prediction of the YneR-PdhA-PdhB complex structure using the power of artificial intelligence suggested that YneR protrudes into the substrate binding site of pyruvate dehydrogenase, thus suggesting a mechanism for inhibition. Indeed, site-directed mutagenesis based on the predicted complex structure verified this mechanism (O'Reilly et al., 2023). This example shows the power of association analyses.

Another interesting group of highly expressed unknown proteins consists of rather small proteins in the range of 47 to 54 amino acids that are encoded in the 5' regions of highly expressed genes and that are associated to an RNA element that is transcribed in the same orientation. Moreover, the occurrence of these proteins is limited to *B. subtilis* and very close relatives in the genus *Bacillus* (see Table 1). All these features are reminiscent of regulatory elements that are involved in mechanisms similar to attenuation. Actually, the BrmB leader peptide of *brmCD* operon shares all properties with respect to protein size, linkage to an RNA element, and occurrence only in *B. subtilis* (Reilman et al., 2014).

## 4 RNA-BINDING PROTEINS

Interactions between proteins and RNA are indispensable for virtually all processes in any living cell. From the biosynthesis of RNA molecules to their degradation, the entire life cycle of RNA is associated to RNA-binding proteins (RBP). RBPs are significantly more conserved across evolution than non-RNA-binding proteins and comprise about 3-11% of the proteome in all domains of life (Gerstberger et al., 2014). The interplay of protein and RNA has traditionally been viewed as the formation of highly dynamic ribonucleoprotein complexes that facilitate multiple regulatory functions including RNA processing, modification, translation or regulation (Dreyfuss et al., 2002; Babitzke et al., 2019). Such 'canonical' RBPs bind RNA specifically via structurally defined RNA-binding domains such as hnRNP homology (KH) domains (Valverde et al., 2008), the S1 domain (Suryanarayana & Subramanian, 1984), RNA recognition motifs (RRM) (Cléry et al., 2008), the RNA-binding domain of transcriptional antiterminators (Stülke, 2002) or DEAD box helicase domains (Linder and Jankowsky, 2011).

The investigation of protein-RNA interactions is typically based on the premise that specific protein domains can bind RNA. This view, however, neglects the ability of RNA adopt highly diverse structures that in principle can be able to bind any molecule. Strikingly, a study aimed at identifying novel RBPs by crosslinking coupled to high resolution mass spectrometry identified proteins crosslinked to RNA that lack any known canonical RNA-binding domain. The discovery of those non-canonical RNA binders such as the phospho-

glycerate kinase, revealed a fundamentally new group of proteins of both known and unknown function to be potential RBPs (Schmidt et al., 2012; Kramer et al., 2014). The existence and identification of those unconventional RBPs harboring unconventional RBDs proposes novel mechanisms of protein-RNA interaction with new biological functions (Hentze et al., 2018). Interestingly, several of the poorly studied but highly expressed proteins of *B. subtilis* that are thought to bind to RNA and/ or the ribosome do not possess known RNA-binding domains. These proteins are YlbN, YlxR, YrzB as well as the poorly studied RNA-binding SpoVG protein (see Table 1, Burke & Portnoy, 2016). It is tempting to speculate that these proteins will define novel RNA-binding domains.

One particularly important open question in the investigation of RNA-binding proteins in *B. subtilis* and other Gram-positive bacteria is related to the base-pairing of small regulatory RNAs with mRNAs. As in other bacteria, small cis- and trans-acting RNA molecules bind to partially complementary mRNAs to control their stability and translation (Ul Haq et al., 2020; Mars et al., 2016). In *E. coli* , the Hfq protein acts as a chaperone that facilitates the base-pairing between only rather short complementary sequences in the two RNA molecules (Vogel and Luisi, 2011). In *B. subtilis* , a shorter Hfq protein is present; however, despite all attempts so far, there is no indication that Hfq also has this function in *B. subtilis* (Dambach et al., 2013; Hämmerle et al., 2014; Rochat et al., 2015). It has recently been proposed that the lack of 5'-3' exonucleolytic RNases in *E. coli* and related bacteria results in the accumulation of mRNA 3' UTR. This may have provided a larger pool of unconstrained RNA sequences that has stimulated the evolution of Hfq function and small RNA (sRNA) regulation. In contrast, the presence of the 5'-3' exoribonuclease RNase J prevents the accumulation of 3' UTR sRNAs in *B. subtilis* . As a consequence, there was no selective pressure on Hfq to evolve to become a mediator of RNA-RNA interactions (Mediati et al., 2021). Since the physical conditions in the *B. subtilis* cell are certainly similar to those in cells of *E. coli* , it seems unlikely that regulation exerted by base-pairing of RNA molecules works without a chaperone. The identification of the protein that takes over this function is a major challenge in the future research on RNA-binding proteins in *B. subtilis* .

Taken together, to understand both basic and complex processes of the cell, extensive research on RNA-protein interactions is required. Even though numerous conventional RBPs have been identified over the past decades, some still remain to be characterized. With the discovery of unconventional RBPs harboring non-canonical RNA-binding domains, the field of understudied RBPs has dramatically increased.

## 5 AMINO ACID TRANSPORT

Although amino acids are the essential building blocks for proteins, there is still a lot to be uncovered, when it comes to the transport of amino acids. Many bacteria including *B. subtilis* can synthesize all 20 proteinogenic amino acids. However, the bacteria do not depend on amino acid synthesis, as they can also directly take up amino acids from their environment. In fact, none of the amino acid biosynthetic pathways is essential in *B. subtilis,* which suggests the presence of transporters for each amino acid. So far, no uptake systems have been identified for asparagine, phenylalanine, glycine, and tyrosine. Several problems make the identification of novel amino acid transporters a challenging task: (i) Multiple transporters exist for most amino acids, (ii) the substrate specificity of many transporters is low, allowing them to take up multiple different amino acids/ metabolites in a rather promiscuous way, and (iii) the affinity of multiple transporters for one amino acid may differ substantially, which results in some transporters only being active under specific conditions. Particularly the presence of multiple transporters for one amino acids (and additionally the existence of biosynthetic pathways) impedes the identification of amino acid transporters since mutants often have no phenotype.

It is important to note that amino acids can also pose a threat to bacteria: histidine inhibits growth of *B. subtilis* in minimal medium (Meißner et al., 2022), and serine and threonine are even toxic to the bacteria under these conditions (Klewing et al., 2020). The toxicity of some amino acids is caused by their high reactivity or by the ability to interfere with other pathways as shown for inhibition of threonine synthesis by serine (de Lorenzo et al., 2015; Mundhada et al., 2017). Amino acid toxicity is even intensified in *B. subtilis* strains lacking the second messenger cyclic di-AMP. This molecule controls potassium homeostasis,

and too high potassium uptake is toxic for the cells (Gundlach et al., 2018). In strains lacking c-di-AMP, potassium is toxic even at low external concentrations in the presence of amino acids. This is caused by the activation of the potassium transporter KtrCD by the common product of amino acid catabolism, glutamate (Krüger et al., 2020). Toxicity and growth inhibition by amino acids can be used to isolate suppressor mutants that have often inactivated the major transport pathway. In this way, the major transporter for serine and glutamate, AimA, was identified (Klewing et al., 2020; Krüger et al., 2021b). Strikingly, the function of this important transporter has remained enigmatic until very recently!

However, it seems that some amino acids do not have one major uptake system, but instead multiple transporters, which contribute to their acquisition. This idea is supported by the fact that suppressors obtained under histidine pressure acquire mutations in the transcriptional repressor AzlB, which allows overexpression of AzlCD, a histidine exporter (Meißner et al., 2022).

Although *B. subtilis* has been studied extensively, there are still several potential amino acid transporters with unknown function (Fig.2). Investigating the remaining uncharacterized transporters could yield valuable insights into overall amino acid metabolism. Moreover, a complete knowledge of amino acid uptake systems would help to understand the requirements that must be met to sustain life in a cell with a minimal genome (Reuß et al., 2016).

Another understudied facet of amino acid transport is the homeostasis of D-amino acids, the enantiomeric counterparts of the proteinogenic L-amino acids. D-amino acids were proposed to act as a bacterial anti-fungal defense mechanism, as they are integrated into fungicidal components like iturin, bacillomycin and mycosubtilin (Stein, 2005). They might also be metabolized by certain bacteria and act as carbon or nitrogen source, which makes them an interesting overall research topic. As far as uptake goes, it is known that certain transporters are able to transport both the D- and L-variant of an amino acid, as is the case for the alanine permease AlaP (Sidiq et al., 2021) The uptake of D-alanine is important, as it is an essential component in the bacterial cell wall. Still, *B. subtilis* is also able to import other D-amino acids such as D-methionine and D-asparagine (Hullo et al., 2004; our unpublished work), which serve no clear purpose within the cell and require further investigation.

## 6 CONTROL OF METABOLIC HOMEOSTASIS

Bacteria are exposed to an ever-changing environment. This forces the cells to constantly sense their surroundings and to adapt to it by regulation of their proteome and metabolome. In this way, bacteria can use their scarce energy resources sparingly and in a targeted manner. Many of the underlying regulatory processes that involve the control of gene expression and enzyme activities are well understood. In contrast, much less is known about the fine-tuning of metabolism which is often achieved by rapidly evolving interactions with regulatory proteins or RNA molecules. As regulation of gene expression is rather slow, direct control of enzyme activities by regulatory interactions seems to be the way of choice for fine-tuning. Indeed, it has been demonstrated that transcription regulation is not sufficient to explain metabolic adaptation (Schilling et al., 2007; Chubukov et al., 2013). One of these examples has recently been identified: the so far unknown B. subtilis protein YneR (see Table 1) was found to interact with two subunits of pyruvate dehydrogenase. Further investigation revealed that this protein inhibits the activity of the pyruvate dehydrogenase (see above; O'Reilly et al., 2023). Similar, the potential control of lipid acquisition by the paralogous proteins YqhY and YloU (Table 1, see above) is an excellent example for the need to study regulatory interactions in metabolism in more details.

A well-established mechanism for fine-tuning a pathway is mediated by regulatory proteins of the PII superfamily. These relatively small proteins can be found in all domains of life and belong to one of the largest family in the group of signal transduction proteins (Forchhammer & Lüddecke, 2016; Forchhammer et al., 2022). Typically, PII proteins bind small molecules to sense the cell's environment and subsequently interact with a wide variety of proteins in the cell. In the first few years after the discovery of the PII protein, it was already shown to be involved in numerous processes of nitrogen anabolism (Ninfa & Jiang, 2005). Through advances in protein structure research, it became apparent that the family of canonical PII proteins has not

only to be expanded to include the PII-like proteins but that these proteins are also involved in many novel regulatory interactions (Forchhammer et al., 2022).

Although the sequence similarity of PII-like proteins is rather low, their trimeric structure is highly conserved and characteristic of this protein family. Excitingly, this has allowed the identification of numerous new sensors that can bind a range of different ligands. For example, the cyanobacterial PII-like signaling protein SbtB was shown to bind numerous adenine nucleotides, including second the messenger molecules cAMP and c-di-AMP. Thus, SbtB not only regulates the bicarbonate transporter SbtA, but also controls the activity of the glycogen-branching enzyme GlgB (Selim et al., 2018; Selim et al., 2021; Fang et al., 2021). Another example is the carboxysome-associated PII-like protein CPII, which binds bicarbonate in addition to ADP/AMP and is thought to regulate carbon metabolism in response to bicarbonate availability (Wheatley et al., 2016). Numerous other PII (-like) proteins, such as the *B. subtilis* DarA protein or CutA from *E. coli* and cyanobacteria, have already been identified and are still waiting to be explored (Gundlach et al., 2015; Selim et al., 2021). Interestingly, DarA binds the essential second messenger molecule c-di-AMP in *B. subtilis* and related gram-positive bacteria (Campeotto et al., 2015; Sureka et al., 2014; Gundlach et al., 2015); however, the function has still not been identified. For CutA, it has long been assumed that it is involved in copper homeostasis; however, a recent study excludes this possibility, thus leaving both the nature of the ligand of CutA and its molecular function an open question (Selim et al., 2021). The fact, that both mutations in *darA* and *cutA* do not yield clear phenotypes supports the idea that these proteins play a role in fine-control of metabolic homeostasis.

Certainly, we are only at the beginning of understanding the mechanisms by which bacterial metabolism can be adapted to subtle changes in the environment. This is important not only to understand the mechanisms that underly the robustness of metabolic networks, but also - due to the ever-growing demand for industrial products, such as drugs, chemicals or vitamins - underlines the importance of developing pathways that are as efficient as possible and thus of our knowledge of metabolic fine-control.

## 7 CONCLUDING REMARKS

Even after decades of research, there are large gaps in our knowledge that impede the comprehensive understanding of life. This applies from the smallest artificial organisms as well as to model organisms such as *M. mycoides* JCVI-syn3A and *B. subtilis*, respectively (Hutchison et al., 2016; Pedreira et al., 2022b). We are now in a position to get novel insights from global unbiased analyses which has already provided (Nicolas et al., 2012; O'Reilly et al., 2023) or will provide at least some initial minimal annotation for many of these proteins. This information is an excellent starting point for further analyses of unknown protein function in *B. subtilis* . The better understanding of these poorly studied proteins as well as of the poorly studied functions will be important to answer key scientific questions such as the identification of a minimal gene set to run a living cell, to develop new strategies in biotechnology, and will also stimulate the research with other bacteria, including important pathogens.
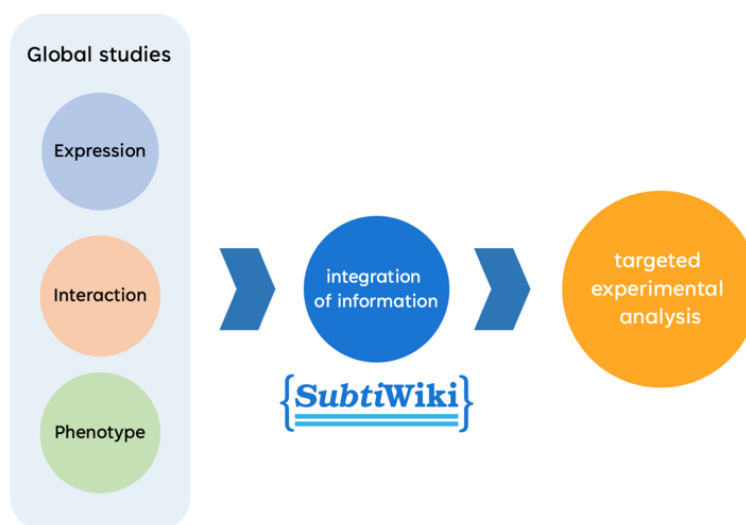
## ACKNOWLEDGEMENTS

## AUTHOR CONTRIBUTIONS

J.S. and C.E. performed data analysis. D.W., J.M., R.W. and J.S. wrote the manuscript. All authors participated in proofreading and corrections.

## GRAPHICAL ABSTRACT

## ABBREVIATED SUMMARY

Bacteria are rather simple forms of life. However, despite extensive research many the functions of many potentially very important proteins have not yet been elucidated even in the model organism *Bacillus subtilis* . We present a workflow from the initial generation of minimal information on such proteins to their functional identification. Moreover, we discuss a set of highly expressed unknown proteins as well as poorly understood biological functions.

## CONFLICT OF INTEREST

Authors declare no conflict of interest.

## REFERENCES

Babitzke, P., Lai, Y.J., Renda, A. & Romeo, T. (2019) Posttranscription initiation control of gene expression mediated by bacterial RNA-binding proteins. *Annual Reviews of Microbiology 73* , 43-67.

Breuer, M., Earnest, T.M., Merryman, C., Wise, K.S., Sun, L., Lynott, M.R. et al. (2019) Essential metabolism for a minimal cell. *Elife8* , e36842.

Burke, T.P. & Portnoy, D.A. (2016) SpoVG is a conserved RNA-binding protein that regulates *Listeria monocytogenes* lysozyme resistance, virulence, and swarming motility. *mBio 7* , e00240.

Campeotto, I., Zhang, Y., Mladenov, M.G., Freemont, P.S. & Gründling, A. (2015) Complex structure and biochemical characterization of the *Staphylococcus aureus* cyclic diadenylate monophosphate (c-di-AMP)-binding protein PstA, the founding member of a new signal transduction protein family. *Journal of Biological Chemistry 290* , 2888-901.

Chihara, K., Gerovac, M., Hör, J., & Vogel, J. (2022) Global profiling of the RNA and protein complexes of *Escherichia coli* by size exclusion chromatography followed by RNA sequencing and mass spectrometry (SEC-seq). *RNA 29* , 123-139.

Chubukov., V., Uhr, M., Le Chat, L., Kleijn, R.J., Jules, M., Link, H. et al. (2013) Transcription regulation is insufficient to explain substrate-induced flux changes in *Bacillus subtilis* .*Molecular Systems Biology 9* ,

709.

Cléry, A., Blatter, M. & Allain, F.H. (2008) RNA recognition motifs: boring? Not quite. *Current Opinion in Structural Biology 18* , 290-298.

Dambach, M., Irnov, I. & Winkler, W.C. (2013) Association of RNAs with *Bacillus subtilis* Hfq. *PLoS one 8* , e55156.

De Jong, L., Roseboom, W. & Kramer, G. (2021) A composite filter for low FDR of protein-protein interactions detected by in vivo cross-linking. *Journal of Proteomics 230* , 103987.

De Lorenzo, V., Sekowska, A, & Danchin, A. (2015) Chemical reactivity drives spatiotemporal organization of bacterial metabolism. *FEMS Microbiology Reviews 39* , 96-119.

Dreyfuss, G., Kim, V. N. & Kataoka, N. (2002) Messenger-RNA-binding proteins and the messages they carry. *Nature Reviews Molecular Cell Biology 3* , 195–205.

Elfmann, C., Zhu, B., Pedreira, T., Hoßbach, B., Lluch-Senar, M., Serrano, L., & Stülke, J. (2022) *Myco* Wiki: functional annotation of the minimal model organism *Mycoplasma pneumoniae* . *Frontiers in Microbiology 13* , 935066.

Fang, S., Huang, X., Zhang, X., Zhang, M., Hao, Y., Guo, H. et al. (2021) Molecular mechanism underlying transport and allosteric inhibition of bicarbonate transporter SbtA. *Proceedings of the National Academy of Sciences of the U. S. A. 118* , e2101632118.

Forchhammer, K. & Lüddecke, J. (2016) Sensory properties of the PII signalling protein family. *FEBS Journal 283* , 425-437.

Forchhammer, K., Selim, K.A. & Huergo LF. (2022) New views on PII signaling: from nitrogen sensing to global metabolic control. *Trends in Microbiology 30* , 722-735.

Galperin, M.Y., Wolf, Y.I., Makarova, K.S., Vera Alvarez, R., Landsman, D. & Koonin, E.V. (2021) COG database update: focus on microbial diversity, model organisms, and widespread pathogens. *Nucleic Acids Research 49* , D274-D281.

Gerstberger, S., Hafner, M. & Tuschl, T. (2014) A census of human RNA-binding proteins. *Nature Reviews Genetics 15* , 829–845.

Gundlach, J., Dickmanns, A., Schröder-Tittmann, K., Neumann, P., Kaesler, J., Kampf, J. et al. (2015) Identification, characterization, and structure analysis of the cyclic di-AMP-binding PII-like signal transduction protein DarA. *Journal of Biological Chemistry 290* , 3069-3080.

Gundlach, J., Commichau, F. M. & Stülke, J. (2018) Perspective of ions and messengers: An intricate link between potassium, glutamate, and cyclic di-AMP. *Current Genetics 64* , 191-195.

Haas, D., Thamm, A.M., Sun, J., Huang, L., Sun, L., Beaudoin, A.W. et al. (2022) Metabolite damage and damage control in a minimal genome. *mBio 13* , e0163022.

Hämmerle, H., Amman, F., Vecerek, B., Stülke, J., Hofacker, I. and Bläsi, U. (2014) Impact of Hfq on the *Bacillus subtilis* transcriptome. *PLoS one 9* , e98661.

Hentze, M. W., Castello, A., Schwarzl, T. & Preiss, T. (2018) A brave new world of RNA-binding proteins. *Nature Reviews Molecular Cell Biology 19* , 327–341.

Hullo, M. F., Auger, S. Dassa, E., Danchin, A. & Martin-Verstraete, I. (2004) The *metNPQ* operon of *Bacillus subtilis* encodes an ABC permease transporting methionine sulfoxide, D- and L-methionine. *Research in Microbiology 155* , 80–86.

Hutchison C.A. 3rd, Chuang R.Y., Noskov V.N., Assad-Garcia, N., Deerinck, T.J., Ellisman, M.H. *et al* . (2016) Design and synthesis of a minimal bacterial genome. *Science 351* , aad6253.

Klewing, A., Koo, B. M., Krüger, L., Poehlein, A., Reuß, D., Daniel, R. et al. (2020) Resistance to serine in *Bacillus subtilis* : identification of the serine transporter YbeC and of a metabolic network that links serine and threonine metabolism. *Environmental Microbiology 22* , 3937–3949.

Koo, B.M., Kritikos, G., Farelli, J.D., Todor, H., Tonk, K., Kimsey, H. et al. (2017) Construction and analysis of two genome-scale libraries for *Bacillus subtilis . Cell Systems 4* , 291-305.

Kramer, K., Sachsenberg, T., Beckmann, B. M., Qamar, S., Boon, K. L., Hentze, M.W. et al. (2014) Photo-cross-linking and high-resolution mass spectrometry for assignment of RNA-binding sites in RNA-binding proteins. *Nature Methods 11* , 1064–1070.

Krüger, L., Herzberg, C., Warneke, R., Poehlein, A., Stautz, J., Weiß, M. et al. (2020) Two ways to convert a low-affinity potassium channel to high affinity: Control of *Bacillus subtilis* KtrCD by glutamate.*Journal of Bacteriology 202* , e00138-20.

Krüger, L., Herzberg, C., Wicke, D., Bähre, H., Heidemann, J. L., Dickmanns, A. et al. (2021a) A meet-up of two second messengers: The c-di-AMP receptor DarB controls (p)ppGpp synthesis in *Bacillus subtilis* . *Nature Communications 12* , 1210.

Krüger, L., Herzberg, C., Rath, H., Pedreira, T., Ischebeck, T., Poehlein, A. et al. (2021b) Essentiality of c-di-AMP in *Bacillus subtilis* : bypassing mutations converge in potassium and glutamate homeostasis. *PLoS Genetics 17* , e1009092.

Krüger, L., Herzberg, C., Wicke, D., Scholz, P., Schmitt, K., Turdiev, A. et al. (2022) Sustained control of pyruvate carboxylase by the essential second messenger cyclic di-AMP in *Bacillus subtilis .mBio 13* , e03602-21.

Kustatscher, G., Colins, T., Gingras, A.C., Guo, T., Hermjakob, H., Ideker, T. et al. (2022a) An open invitation to the Understudied Proteins Initiative. *Nature Biotechnology 40* , 815-817.

Kustatscher, G., Colins, T., Gingras, A.C., Guo, T., Hermjakob, H., Ideker, T. et al. (2022b) Understudied Proteins: opportunities and challenges for functional proteomics. Nature *Methods 19* , 774-779.

Lamm-Schmidt, V., Fuchs, M., Sulzer, J., Gerovac, M., Hör, J., Dersch, P. et al. (2021) Grad-seq identifies KhpB as a global RNA-binding protein in *Clostrioides difficile . microLife 2* , uqab004.

Linder, P. & Jankowsky, E. (2011) From unwinding to clamping – the DEAD box RNA helicase family. *Nature Reviews Molecular Cell Biology 12* , 505-516.

Link, H., Kochanowski, K., & Sauer, U. (2013) Systematic identification of allosteric protein-metabolite interactions that control enzyme activity *in vivo* . *Nature Biotechnol* ogy *31* , 357-361.

Mars, R.A., Nicolas, P., Denham, E.L. & van Dijl, J.M. (2016) Regulatory RNAs in *Bacillus subtilis* : a gram-positive perspective on bacterial RNA-mediated control of gene expression. *Microbiology and Molecular Biology Reviews 80* , 1029-1057.

Mediati, D.G., Lalaouna, D. and Tree, J.J. (2021) Burning the candle at both ends: have exoribonucleases driven divergence of regulatory RNA mechanisms in bacteria? *mBio 12* , e0104121.

Meißner, J., Schramm, T., Hoßbach, B., Stark, K., Link, H. & Stülke, J. (2022) How to deal with toxic amino acids: the bipartite AzlCD complex exports histidine in *Bacillus subtilis* . *Journal of Bacteriology 204* , e00353-22.

Mundhada, H., Seoane, J.M., Schneider, K., Koza, A., Christensen, H.B., Klein, T. et al. (2017) Increased production of L-serine in*Escherichia coli* through adapative laboratory evolution.*Metabolic Engineering 39* , 141-150.

Nicolas, P., Mäder, U., Dervyn, E., Rochat, T., Leduc, A., Pigeonneau, N. et al. (2012) Condition-dependent transcriptome reveals high-level regulatory architecture in *Bacillus subtilis* . *Science 335* , 1103-1106.

Niehaus, T.D., Elbadawi-Sidhu, M., de Crécy-Lagard, V. Fiehn, O. & Hanson, A.D. (2017) Discovery of a widespread prokaryotic 5-oxoprolinase that was hiding in plain sight. *Journal of Biological Chemistry 292* , 16360-16367.

Ninfa, A.J. & Jiang, P. (2005) PII signal transduction proteins: sensors of alpha-ketoglutarate that regulate nitrogen metabolism. *Current Opinion in Microbiology 8* , 168-173.

O'Reilly, F.J., Xue, L., Graziadei, A., Sinn, L., Lenz, S., Tegunov, D. et al. (2020) In-cell architecture of an actively transcribing-translating expressome. *Science 369* , 554-557.

O'Reilly, F.J., Graziadei, A., Forbrig, C., Bremenkamp, R., Charles, K., Lenz, S. et al. (2023) Protein complexes in cells by AI-assisted structural proteomics. doi.org/10.1101/2022.07.26.501605.

Pedreira, T., Elfmann, C., Singh, N., & Stülke, J. (2022a) *Syn* Wiki: Functional annotation of the first artificial organism *Mycoplasma mycoides* JCVI-syn3A. *Protein Science 31* , 54-62.

Pedreira, T., Elfmann, C., and Stülke, J. (2022b) The current state of *Subti* Wiki, the database for the model organism *Bacillus subtilis* . *Nucleic Acids Research 50* , D875-D882.

Peters, J.M., Colavin, A., Shi, H., Czarny, T.L., Larson, M.H., Wong, S. et al. (2016) A comprehensive, CRISPR-based functional analysis of essential genes in bacteria. *Cell 165* , 1493-1506.

Reilman, E., Mars, R.A., van Dijl, J.M. & Denham, E.L. (2014) The multidrug ABC transporter BmrC/BmrD of *Bacillus subtilis* is regulated via a ribosome-mediated transcriptional attenuation mechanism. *Nucleic Acids Research 42* , 11393-11407.

Reuß, D. R., Commichau, F. M., Gundlach, J., Zhu, B. & Stülke, J. (2016) The blueprint of a minimal cell: *MiniBacillus* . *Microbiology and Molecular Biology Reviews 80* , 955-987.

Reuß, D. R., Altenbuchner, J., Mäder, U., Rath, H., Ischebeck, T., Sappa, P. K. et al. (2017) Large-scale reduction of the *Bacillus subtilis* genome: consequences for the transcriptional network, resource allocation, and metabolism. *Genome Research* 27, 289-299.

Rochat, T., Delumeau, O., Figuera-Bossi, N., Noirot, P., Bossi, L., Dervyn, E. & Bouloc, P. (2015) Tracking the elusive function of *Bacillus subtilis* Hfq. *PLoS one 10* , e0124977.

Roney, I.J. & Rudner, D.Z. (2023) Two broadly conserved families of polyprenyl-phosphate transporters. Nature in press. doi.org/10.1038/s41586-022-05587-z.

Sachla, A. & Helmann, J.D. (2019) A bacterial checkpoint protein for ribosome assembly moonlights as an essential metabolite-proofreading enzyme. *Nature Communications 10* , 1526.

Saha, C.K., Pires R.S., Brolin, H., Delannoy, M. & Atkinson, G.C. (2021) FLaGs and webFLaGs: discovering novel biology through the analysis of gene neighbourhood conservation. *Bioinformatics 37* , 1312-1314.

Schilling, O., Frick, O., Herzberg, C., Ehrenreich, A., Heinzle, E., Wittmann, C. & Stülke, J. (2007) Transcriptional and metabolic responses of Bacillus subtilis to the availability of organic acids: transcription regulation is important but not sufficient to account for metabolic adaptation. *Applied and Environmental Microbiology 73* , 499-507.

Schmidt, C., Kramer, K. & Urlaub, H. (2012) Investigation of protein-RNA interactions by mass spectrometry - Techniques and applications. *Journal of Proteomics 75* , 3478–3494.

Selim, K.A., Haase, F., Hartmann, M.D., Hagemann, M. & Forchhammer, K. (2018) $P_{II}$-like signaling protein SbtB links cAMP sensing with cyanobacterial inorganic carbon response. *Proceedings of the National Academy of Sciences of the U. S. A. 115* , E4861-E4869.

Selim, K. A., Haffner, M., Burkhardt, M., Mantovani, O., Neumann, N., Albrecht, R. et al. (2021) Diurnal oscillation of cyanobacteria requires perception of second messenger signaling molecule c-di-AMP by the carbon-control protein SbtB. Science Advances 7, eabk0568.

Senges, C.H.R., Stepanek, J.J., Wenzel, M., Raatschen, N., Ay, Ü., Märtens, Y. et al. (2021) Comparison of proteomic responses as global approach to antibiotic mechanism of action elucidation. *Antimicrobial Agents and Chemotherapy 64* , e01373-20.

Senges, C.H.R., Warmuth, H.L., Vázquez-Hernández, M., Uzun, H.D., Sagurna, L., Dietze, P. et al. (2022) Effects of 4-Br-A23187 on *Bacillus subtilis* cells and unilamellar vesicles reveal it to be a potent copper ionophore. *Proteomics 22* , e2200061.

Sidiq, K. R., Chow, M. W., Zhao, Z. & Daniel, R. A. (2021) Alanine metabolism in *Bacillus subtilis* . *Molecular Microbiology 115* , 739–757.

Sklarczyk, D., Kirsch, R., Koutrouli, M., Nastou, K., Mehryary, F., Hachilif, R. et al. (2023) The STRING database in 2023: protein-protein association networks and functional enrichment analyses for any sequenced genome of interest. *Nucleic Acids Research 51* , D638-D646.

Stein, T. (2005) *Bacillus subtilis* antibiotics: structures, syntheses and specific functions. *Molecular Microbiology 56* , 845–857.

Stülke, J. (2002) Control of transcription termination in bacteria by RNA-binding proteins by RNA-binding proteins that modulate RNA structures. *Archives of Microbiology 177* , 433-440.

Sureka, K., Choi, P.H., Precit, M., Delince, M., Pensinger, D.A., Huynh, T.N. et al. (2014) The cyclic dinucleotide c-di-AMP is an allosteric regulator of metabolic enzyme function. *Cell 158* , 1389-1401.

Suryanarayana, T. & Subramanian, A.R. (1984) Function of the homologous sequences in nucleic acid binding domain of ribosomal protein S1. *Biochemistry 23* , 1047-1051.

Tödter, D., Gunka, K. & Stülke, J. (2017) The highly conserved Asp23 family protein YqhY plays a role in lipid biosynthesis in *Bacillus subtilis* . *Frontiers in Microbiology 8* , 883.

Ul Haq, I., Müller, P. &. Brantl, S. (2020). Intermolecular communication in *Bacillus subtilis* : RNA-RNA, RNA-protein and small protein-protein interactions. *Frontiers in Molecular Biosciences 7* , 178.

UniProt Consortium (2023) UniProt: the universal protein knowledgebase in 2023. *Nucleic Acids Research 51* , D523-DD531.

Valverde, R., Edwards, L. & Regan, L. (2008) Structure and function of KH domains. *FEBS Journal 275* , 2712–2726.

Vogel, J. & Luisi, B.F. (2011) Hfq and its constellation of RNA. *Nature Reviews Microbiology 9* , 578-589.

Wheatley, N.M., Eden, K.D., Ngo, J., Rosinski, J.S., Sawaya, M.R., Cascio, D. et al. (2016) A PII-like protein regulated by bicarbonate: Structural and biochemical studies of the carboxysome-associated CPII protein. *Journal of Molecular Biology 428* , 4013-4030.

Yang, J., Suzuki, M. & McCarty, D.R. (2016) Essential role of conserved DUF177a protein in plastid 23S rRNA accumulation and plant embryogenesis. *Journal of Experimental Botany 67* , 5447-5460.

Table 1. A set of highly expressed proteins of unknown function in *B. subtilis* [1]

| Protein | Average expression | Occurence | Available information | Association to known genes/ proteins |
|---|---|---|---|---|
| Potential RNA/ ribosome-binding proteins | Potential RNA/ ribosome-binding proteins | Potential RNA/ ribosome-binding proteins | Potential RNA/ ribosome-binding proteins | Potential RNA/ ribosome-binding proteins |

| Protein | Average expression | Occurence | Available information | Association to known genes/ proteins |
|---|---|---|---|---|
| YqeY | 16.09 | Majority of bacteria | may have tRNA amino acid amidase activity (COG) | *rpsU* (operon) |
| YtpR | 13.08 | Bacilli, some Mollicutes | tRNA-binding domain, interacts with GatB (IAX, O'Reilly), may participate in GatABC activity | GatB (O'Reilly et al., 2023) |
| YabR | 14.31 | Many Firmicutes and few other bacteria | RNA-binding S1 domain, ribosome-interacting protein (O'Reilly), paralog of **YugI** | *rcqP*, *divIC* (operon) RcqU (O'Reilly et al., 2023) |
| YugI | 14.13 | Many Firmicutes and few other bacteria | RNA-binding S1 domain, ribosome-interacting protein (O'Reilly), paralog of **YabR** | Pgk, EzrA (O'Reilly et al., 2023) |
| YrzB | 13.30 | Many Firmicutes | Associated to the ribosome (IAX, O'Reilly) | RplR, RplT, RpmF, RpsN, RpsL (O'Reilly et al., 2023) *reoM* (operon) |
| YlxR | 13.90 | Widespread in bacteria, but not in beta- and gamma-Proteobacteria | May bind to RNA (by similarity, Lamm-Schmidt et al., 2021) | *nusA*, *rplGA*, *infB*, *rbfA* (operon) |
| KhpA | 13.74 | Widespread in bacteria, but not in alpha-, beta-, gamma-Proteobacteria | KH domain, putative RNA-binding protein | *ffh*, *rpsP*, *rimM*, *trmD*, *rplS* (operon) |
| YlbN | 15.42 | Most bacteria and plant chloroplasts | essential in plants, required for 23S rRNA accumulation | *rpmF* (operon) RpsJ (O'Reilly et al., 2023) |
| **Potentially involved in the control of metabolism** | **Potentially involved in the control of metabolism** | **Potentially involved in the control of metabolism** | **Potentially involved in the control of metabolism** | **Potentially involved in the control of metabolism** |

| Protein | Average expression | Occurence | Available information | Association to known genes/ proteins |
|---|---|---|---|---|
| YjlC | 14.52 | Some bacteria and archaea | operon with *ndh*, and interaction with Ndh, may control Ndh activity | Ndh (operon, O'Reilly et al., 2023) |
| YlaN | 14.00 | Many Bacilli | Essential in the absence of iron, putative Fur effector | Fur (O'Reilly et al., 2023) |
| YqhY | 13.99 | Many bacteria (not in proteobacteria) | Asp23 family, paralog of **YloU**, conserved gene cluster with *accBC*, quasi-essential, may control fatty acid synthesis | *accB*, *accC* (operon) |
| YloU | 13.18 | Many bacteria (not in proteobacteria) | Asp23 family, paralog of **YqhY**, conserved gene cluster with *fakA*, may control fatty acid phosphorylation | *fakA* (operon) |
| YneR | 13.12 | Bacilli | Interaction with PdhA and PdhB, inbits PDH activity (IAX, O'Reilly), renamed PdhI | PhdA, PdhB (O'Reilly et al., 2023) *plsY* (operon) |
| **Potential regulatory elements** | **Potential regulatory elements** | **Potential regulatory elements** | **Potential regulatory elements** | **Potential regulatory elements** |
| YtzK | 14.32 | *B. subtilis* | May be regulatory element for *tyrS* expression | *acsA*, *tyrS* (operon) |
| YpzE | 13.12 | *B. subtilis* | Upstream of *ribU*, may control expression of *ribU* | *ribU* (operon) |
| YwzH | 13.76 | *Bacillus* sp. | First gene of the *dlt* operon for teichoic acid biosynthesis, microprotein, may control expression of the operon | *dltABCDE* (operon) |
| YqzL | 13.73 | *Bacillus* sp. | Microprotein | *cdd*, *era*, *recO* (operon) |

15

| Protein | Average expression | Occurence | Available information | Association to known genes/ proteins |
|---|---|---|---|---|
| Other highly expressed unknown proteins | Other highly expressed unknown proteins | Other highly expressed unknown proteins | Other highly expressed unknown proteins | Other highly expressed unknown proteins |
| YkaA | 13.50 | Many bacteria and archaea | Putative phosphate transport regulator (COG) | *pit* (operon) |
| YjcF | 13.49 | Widespread in bacteria & euryarchaeota | Predicted N-acyltransferase, GNAT family, operon with yjcH | Pta, RplW (O'Reilly et al., 2023) |
| YjcH | 13.43 | Widespread in bacteria | Operon with yjcF | |
| YhaH | 13.35 | Widespread in bacteria, but not in alpha-, beta-, gamma-Proteobacteria | Membrane protein, paralog of **YtxH** | AccA (O'Reilly et al., 2023) |
| YtxH | 13.99 | Widespread in bacteria, but not in alpha-, beta-, gamma-Proteobacteria | Paralog of **YhaH** | *brxC* (operon) |
| YqgA | 15.22 | Closely related *Bacillus* species | attached to cell wall, division sites | |
| YybN | 14.66 | *B. subtilis* | | |
| YeeI | 13.32 | Nearly all bacteria | Putative regulatory protein (COG) | |
| YydD | 13.24 | Very few bacteria | Low expression during sporulation | |
| YhjA | 13.80 | Some *Bacillus* sp. | | |
| YpjP | 13.43 | *Bacillus* sp. | Low expression during sporulation | *thyB*, *dfrA* (operon) |
| YkuJ | 14.01 | Many Bacilli | Operon with ykuK, *abbB*, *darB* and *ccpC* | NusA (O'Reilly et al., 2023) *abbB*, *darB, ccpC* (operon) |
| YkuK | 13.22 | Many Bacilli | Operon with ykuJ, *abbB*, *darB* and *ccpC* | *abbB, darB, ccpC* (operon) |
| YocA | 13.10 | Most bacteria | Membrane protein, glycoside hydrolase family 23 (UinProt) | |

| Protein | Average expression | Occurence | Available information | Association to known genes/ proteins |
|---|---|---|---|---|
| YqkB | 13.06 | Many Bacilli | Predicted Fe-S cluster biosynthesis protein (COG) | |
| YneF | 14.51 | Most Bacilli and all Mollicutes | Membrane protein, essential | *sirA* (operon) |
| YubF | 13.02 | *Bacillus* sp. | Membrane protein | *lytG* (operon) |
| YfhC | 13.03 | Most bacteria and archaea | Nitroreductase (COG) | |
| YqxD | 14.10 | Many Firmicutes and Proteobacteria | | *dnaG*, *sigA* (operon) |
| YusG | 14.09 | *Bacillus* sp. | | *gcvH* (operon) |
| YaaK | 13.95 | Widespread in all bacteria | Putative DNA-binding protein | *dnaX*, recR, *bofA* (operon) |
| YhjA | 13.80 | Some *Bacillus* sp. | | |
| YlmG | 13.88 | Many bacteria | Membrane protein, cytochrome maturation (COG) | sepF (operon) |
| YaaR | 13.55 | Firmicutes, Spirochaetes, Thermotoga | | *tmk*, *holB*, *trmN6* (operon) |
| YbzG | 15.78 | Very few Firmicutes | part of ribosomal protein/ translation factor/ *rpoA* operon | |

Arbitrary expression levels as introduced in Nicolas et al. (2012) were used. The number corresponds to logarithmical presentation of the experimentally determined expression. To select this set of proteins, we used two conditions: First, the average expression level under 104 conditions should be above 13. Second, the gene should be highly expressed under all conditions. Exceptions were allowed for low expression during sporulation. Proteins were regarded unknown if no function had be identified prior to a large scale *in vivo* interaction study (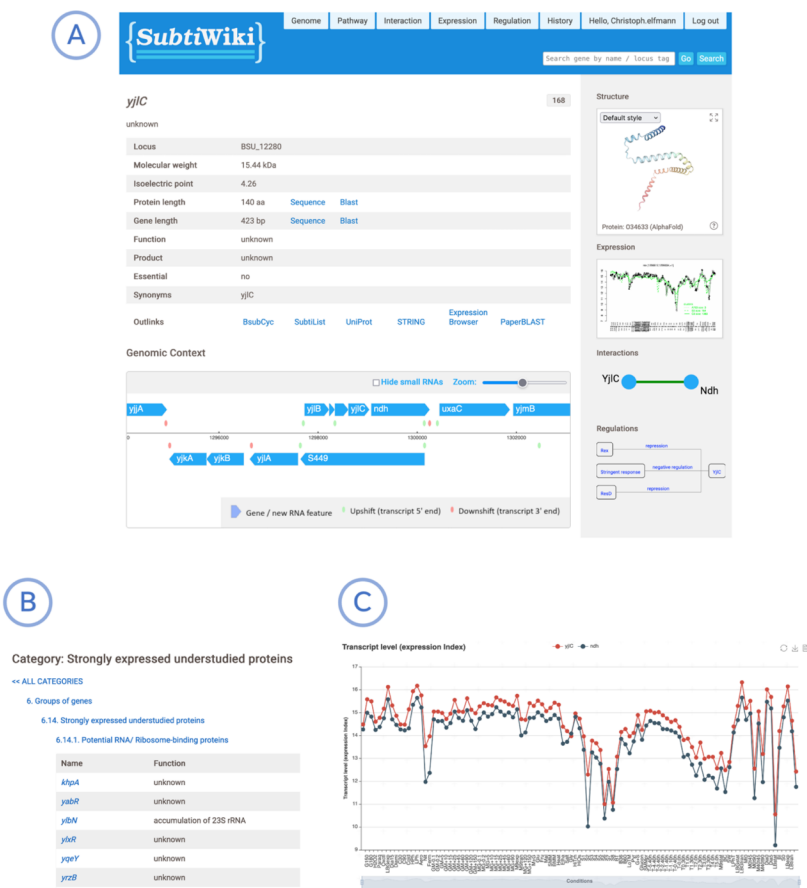O'Reilly et al., 2023). The list is also accessible online (http://www.subtiwiki.uni-goettingen.de/v4/category?id=SW.6.1).

**Figure Legends**

**Figure 1. Integration of information on *Bacillus subtilis* genes and proteins in the database *Subti*Wiki (http://subtiwiki.uni-goettingen.de).**

1. Basic information on the highly expressed unknown protein YjlC. The gene-specific page includes all available information on the protein as well as graphical interactive elements for the genomic content, the structure of the protein, the expression of the gene under more than 100 conditions, protein-protein interactions, and regulation of the gene. Note that the *yjlC* and *ndh* genes form an operon, and the

encoded proteins interact physically (see also Table 1) indicating a functional link between the two proteins.

2. The *Subti* Wiki category for highly expressed understudied proteins give immediate access to the set of 41 proteins, sorted according to potential functions (see Table 1).

3. The Expression Browser of SubtiWiki allows a direct comparison of the expression levels of two or more genes (here *yjlC* and *ndh* ). Similar expression patterns support the idea of a functional link between two genes/proteins.

**Figure 2 . Overview of transporter proteins with unknown function and amino acids with no known transporter** . While some of these proteins are only expressed during sporulation, others are moderately expressed under standard growth conditions. For five of twenty proteinogenic amino acids no transporter responsible for the uptake of the particular amino acid could be identified. The existence of these transporters is still guaranteed, as mutants of the synthesis genes for these five amino acids are viable in complex medium.



Wicke et al., Fig. 1

**Hosted file**

`image3.emf` available at https://authorea.com/users/582320/articles/622445-understudied-proteins-and-understudied-functions-in-the-model-bacterium-bacillus-subtilis-a-major-

challenge-in-current-research

Wicke et al., Fig. 2