1    **Data and Text Mining**

2    **TransporterPAL: An integrative database Transporter Prediction ALgorithm**

3    Jane Dannow Dyekjær[1], Alexander Kruhøffer Bloch Andersen[1], Joel August Vest Madsen[2], Jens

4    Preben Morth[3], Lars Juhl Jensen[4], Irina Borodina[1,*]

5    [1]Novo Nordisk Foundation Center for Biosustainability, Technical University of Denmark,

6    Building 220, Kemitorvet, DK-2800, Kgs. Lyngby, Denmark

7    [2]VesTech, A N Hansens Alle 31a, DK-2900 Hellerup. Denmark

8    [3]DTU Bioengineering, Søltofts Plads, Building 221, DK-2800 Kgs. Lyngby, Denmark

9    [4]Novo Nordisk Foundation Center for Protein Research, Faculty of Health and Medical Sciences,

10    University of Copenhagen, Blegdamsvej 3B, DK-2200 Copenhagen N, Denmark

11    *To whom correspondence should be addressed.

12    **Abstract**

13    **Motivation:** Natural products are used as drugs, cosmetic ingredients, pigments, flavors, and

14    agricultural products. The compounds are retrievable as extracts from natural sources, but the

15    yields are often low, and the final product may contain various impurities. These challenges can

16    be solved by expressing the biosynthetic pathway in microbial cell factories and ensuring

17    product secretion from the cell by using an appropriate transporter. However, insufficient

18    knowledge of transporters for specific compounds often obstructs efficient secretion of the

19    natural product. Therefore, our goal was to develop an algorithm that predicts transporters for a

20    given compound using available public data.

21    **Results:** The web application TransporterPAL predicts suitable transporters for compounds by

22    interconnecting data for biosynthetic genes and their interactions with transporters. The web

23    application queries the STITCH, STRING, and UniProtKB databases via their respective APIs

24  and returns a set of potential transporters based on a compound and, optionally, the organism as

25  input. For a test set of 61 transporter systems, each containing one or more transporters, a total of

26  90 unique transporters with a known substrate, we could retrieve 45% of the transporters. To our

27  knowledge, this is the first bioinformatics tool for predicting transporter candidates for a given

28  molecule.

29  **Availability:** https://transporterpal.com

30  **Contact:** irbo@biosustain.dtu.dk

31  **Supplementary information:** Supplementary data are available at Bioinformatics online.

32  **1   Introduction**

33  The demand for natural products is increasing in the pharmaceutical and biotechnological

34  industries. Due to the complexity of these molecules, the chemical synthesis of natural products is

35  usually challenging, if not impossible. Alternatively, natural products can be produced cheaper

36  and more sustainably by incorporating heterologous pathway genes into microbial cell factories.

37  Secretion of the synthesized molecules from microbial cells is vital to alleviate cellular toxicity,

38  prevent product degradation, and allow for high titer and purity (van der Hoek and Borodina,

39  2020). Often, a specific transmembrane transporter protein is required to engineer the secretion of

40  small molecules. However, transmembrane proteins are challenging to study, and many

41  transporters are thus not well characterized. As a result, the knowledge of appropriate transporters

42  is scarce, but the exploitation of data from several independent databases can provide clues about

43  potential transporters.

44  We hypothesized that if we can extract pathway genes for biosynthesis of a compound and other

45  proteins interacting with the compound, we can identify a suitable transporter among their

46  interaction partners in a functional association network. Furthermore, assuming the transporter is

47 present within this expanded set of proteins, we can use text mining to differentiate transporters

48 from other pathway proteins based on their functional description. Here, we present a web resource

49 that interacts with public databases to predict potential transporters automatically for a given

50 compound produced in nature, irrespective of protein annotation.

51 **2  Implementation**

52 STITCH (Szklarczyk *et al.*, 2016)  is a database containing information about chemical–protein

53 interactions from many sources, including biological pathway databases, automatic text mining of

54 biomedical literature, and experimental data repositories. The STRING database (Szklarczyk *et*

55 *al.*, 2021) similarly provides physical and functional protein-protein interactions integrated from

56 various sources. Finally, UniProtKB (Bateman *et al.*, 2021) provides functional annotations of

57 proteins, including Gene Ontology (GO) annotations related to transport activity.

58 Our workflow consists of the following steps, starting from a chemical and, optionally, an

59 organism of interest: (i) use chemical–protein interactions from STITCH to obtain an initial set of

60 proteins (Szklarczyk *et al.*, 2016), (ii) expand this set of proteins with their interaction partners

61 using functional associations from STRING (Szklarczyk *et al.*, 2021), (iii) select the transport-

62 related proteins from this set based on UniProtKB GO annotations (Figure 1). To combine the

63 information from all these databases, we access them via their respective REST APIs. Specifically,

64 we use the *interactors* API method of STITCH to obtain interacting proteins with the input

65 chemical, then the *interaction_partners* API method of STRING to expand the set of proteins. We

66 convert the STRING identifiers to UniProt IDs using the UniProt Retrieve/ID mapping service and

67 subsequently retrieve the UniProtKB record for each ID via the UniProt API. Each record contains

68 information about the protein, and as we search specifically for transporters, we filter the records

69 based on keywords related to transport activity. We have used UniProt GO annotations for cellular

70 components, molecular functions, and biological processes to generate a list of transporter-related

71 keywords. After removing protein entries unrelated to transport, the UniProtKB accession number,

72 protein name, and organism name are retrieved from UniProtKB. The interaction between the

73 compound and each putative transporter is scored using the combined score from the network

74 STITCH API, as described in (Szklarczyk *et al.*, 2016).

75 The web application is based on Python3 and uses STITCH, STRING, and UniProt APIs. We use

76 the python packages requests, csv, json and sys. Furthermore, asyncio and aiohttp are used to make

77 asynchronous API calls. The web application has been developed using an Express NodeJS web

78 application framework in combination with the basic web application building blocks, HTML,

79 CSS, and Javascript.

80 **3  Availability**

81

82 On the online server https://transporterpal.com, the user enters a natural compound and, optionally,

83 an organism. A link to available organisms is provided on the website. When the web application

84 has finished, it sends an email containing a CSV file and a FASTA file containing the amino acid

85 sequences of the identified transporters.

86

87 **4  Current benchmarking**

88 Although the algorithm works with any organism, we tested the algorithm on a data set of bacterial

89 transporters extracted from the Transporter Classification Database (Saier *et al.*, 2021), a database

90 containing experimentally verified transporters, to quantify the number of transporters we can

91 retrieve by using the algorithm. We chose transporters of degradation pathways and secondary

92    metabolite biosynthesis produced by specific bacteria and corroborated their pathway in MetaCyc

93    (Caspi *et al.*, 2020). We selected transporters from 61 transporter classification IDs, each

94    containing one or more transporters, in total 95 transporters (Supplementary data). A few of the

95    transporters mediate transport of more than one compound, which makes 90 of the transporters

96    unique. Of the unique transporters, 35 were not present in the databases; For 20 compounds,

97    transporters related to a specific organism were unavailable in STITCH, primarily due to

98    unrepresented species. Additional nine transporters were not present in the STRING database.

99    Furthermore, six transporters did not have a link between the STRING and UniProtKB databases.

100    From the remaining data set of 55 proteins, 14 transporters were not found either because they

101    were not among the interaction partner proteins or because they did not contain the transporter-

102    related GO-annotation keywords. Ultimately, 41 proteins were correctly retrieved, corresponding

103    to 46% of the 90 unique transporters. The predictive power of TransporterPAL depends on

104    available data, highlighting the need for more experimental research on transporter function

105    characterization.

106

## 5   Conclusion

108    We have developed a web application to predict potential transporters based on the pathway genes

109    and interacting proteins of a given compound and organism. With more research on transporters

110    and the addition of knowledge into the STITCH, STRING, and UniProtKB databases, we expect

111    to increase the success rate of retrieving suitable proteins with transport activity for a specific

112    compound. While we focus on transporters, this method is easily extendable to study other

113    pathway proteins by changing the filters in the GO annotation to keywords suitable for a particular

114     protein family class. To our knowledge, this is the first tool able to link potential transporters to a

115     given natural product.

120

121     **Conflict of Interest***:*

122     None declared.

123     **References**

124     Bateman,A. *et al.* (2021) UniProt: The universal protein knowledgebase in 2021. *Nucleic Acids*

125     *Res.*, **49**, D480–D489.

126     Caspi,R. *et al.* (2020) The MetaCyc database of metabolic pathways and enzymes-a 2019 update.

127     *Nucleic Acids Res.*, **48**, D455–D453.

128     van der Hoek,S.A. and Borodina,I. (2020) Transporter engineering in microbial cell factories: the

129     ins, the outs, and the in-betweens. *Curr. Opin. Biotechnol.*, **66**, 186–194.
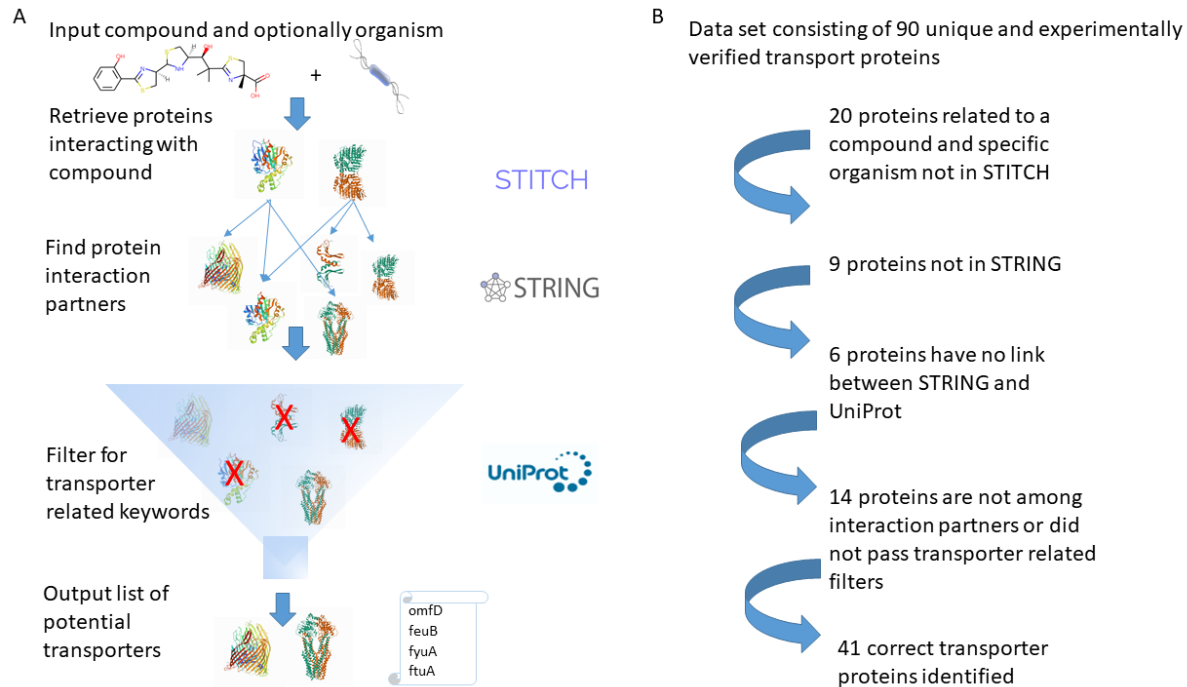
130     Saier,M.H. *et al.* (2021) The transporter classification database (TCDB): 2021 update. *Nucleic*

131     *Acids Res.*, **49**, D461–D467.

132     Szklarczyk,D. *et al.* (2016) STITCH 5: augmenting protein-chemical interaction networks with

133       tissue and  affinity data. *Nucleic Acids Res.*, **44**, D380-D384.

134    Szklarczyk,D. *et al.* (2021) The STRING database in 2021: customizable protein-protein

135       networks, and functional  characterization of user-uploaded gene/measurement sets. *Nucleic*

136       *Acids Res.*, **49**, D605–D612.

137

138

139    Figure 1. TransporterPAL workflow for transporter identification. Panel A. We use the STITCH

140    database to extract pathway genes and proteins interacting with the compound for an organism.

141    This set of proteins is expanded by their interaction partners using STRING. The list contains

142    non-transporter-related proteins, and we exclude these by filtering for transporter-related GO

143    annotations using UniProt to obtain a final list of potential transporters for the compound. Panel

144    B. 90 unique, experimentally verified transporter proteins were used in the test set. Twenty

145    proteins related to a compound and specific organism were unavailable in STITCH. Nine

146    proteins were not present in the STRING, and six transporters did not have a link between the

147    STRING and UniProt databases. In the remaining data set of 55 transporters, 14 were not found,

148    resulting in 41 correctly identified transporter proteins.