

1 **BtToxin_Digger: a comprehensive and high-throughput pipeline for**
2 **mining toxin protein genes from *Bacillus thuringiensis***

3

4

5 Hualin Liu¹, Jinshui Zheng^{1,2*}, Yun Yu¹, Weixing Ye¹, Donghai Peng¹,

6 Ming Sun^{1*}

7

8

9 ¹State Key Laboratory of Agricultural Microbiology, Huazhong Agricultural

10 University, Wuhan, 430070, China

11 ²Hubei Key Laboratory of Agricultural Bioinformatics, Huazhong Agricultural

12 University, Wuhan, 430070, China.

13 * corresponding author

14

15

16 *Corresponding author: E-mail: jszheng@mail.hzau.edu.cn;

17 m98sun@mail.hzau.edu.cn.

18 **Summary:** *Bacillus thuringiensis* (Bt) which is a spore-forming gram-positive
19 bacterium, has been used as the most successful microbial pesticide for decades. Its
20 toxin genes (*cry*) have been successfully used for the development of GM crops against
21 pests. We have previously developed a web-based insecticidal gene mining tool
22 BtToxin_scanner, which has been proved to be the most important method for mining
23 *cry* genes from Bt genome sequences. To facilitate efficiently mining major toxin genes
24 and novel virulence factors from large-scale Bt genomic data, we re-design this tool
25 with a new workflow. Here we present BtToxin_Digger, a comprehensive, high-
26 throughput, and easy-to-use Bt toxin mining tool. It runs fast and can get rich, accurate,
27 and useful results for downstream analysis and experiment designs. Moreover, it can
28 also be used to mine other targeting genes from large-scale genome and metagenome
29 data with the addition of other query sequences.

30 **Availability and Implementation:** The BtToxin_Digger codes and instructions are
31 freely available at https://github.com/BMBGenomics/BtToxin_Digger.

32 **Contact:** jszheng@mail.hzau.edu.cn; m98sun@mail.hzau.edu.cn.

33 **1 Introduction**

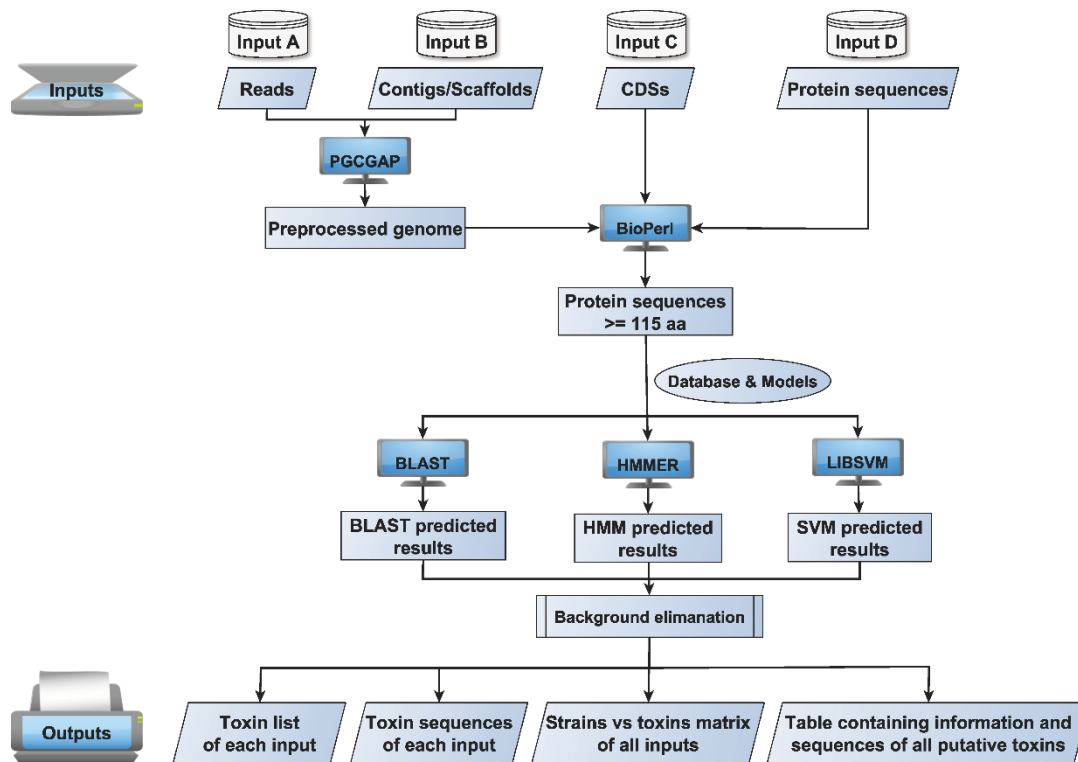
34 The toxins produced by *Bacillus thuringiensis* (Bt) have insecticidal activity against
35 many agricultural and forestry pests, so they are widely used in the development of
36 biopesticides and GM insect-resistant crops. Bt products represent more than 60% of
37 the biopesticide market (Siegwart *et al.*, 2015). Crystal protein (Cry) produced by Bt as
38 the major toxin can kill insects from many orders including Lepidoptera, Diptera, and
39 Coleoptera, etc. The *cry* gene is one of the most important genes used for the
40 development of genetically modified (GM) crops targeting insect pests. From 1996 to

41 2016, the planting of Bt maize and cotton had delivered \$50.6 billion and \$54 billion
42 of extra farm income, respectively (G Brookes and Barfoot, 2018). Due to the
43 importance of Bt toxins, many researchers and companies have been working on the
44 discovery of new toxin genes (Sanahuja *et al.*, 2011). Other toxins with insecticidal
45 activity produced by Bt include Cyt (Cytotoxic toxin protein) and Vip (Vegetative
46 insecticidal protein), etc (Palma *et al.*, 2014). Previously, we developed an on-line tool
47 BtToxin_scanner to predict Crys encoding genes from Bt genome sequences (Ye *et al.*,
48 2012). It can handle several assembled genomes every time and provide useful
49 comparative results between the predicted toxin and with known ones. During the past 7
50 years, it was widely used by researchers worldwide (Méric *et al.*, 2018; Prado *et al.*,
51 2014; Ruan *et al.*, 2015; Zheng *et al.*, 2017). Here we re-designed the previous tool to
52 provide a novel, high-throughput, and local software BtToxin_Digger which can be
53 directly used to handle large-scale genomic and metagenomic data to predict all kinds
54 of putative toxin genes. It also generates comprehensive and readable results to
55 facilitate the downstream sequence analysis or experiment design (Figure 1).

56 **2 Methods**

57 The tool accepts multiple forms of input data including Reads (pair-end reads, long-
58 reads, or hybrid-reads), genomic or metagenomic assemblies, coding sequences (CDSs),
59 and protein sequences. PGCGAP (Liu *et al.*, 2020) was used for genome assembly and
60 pretreatment. ORFs finding and translation are performed by BioPerl (Stajich *et al.*,
61 2002). All protein sequences with a length above 115-aa are searched against the
62 database and trained models by BLAST (Camacho *et al.*, 2009), HMMER (Eddy, 2011),

63 and LIBSVM (Chang and Lin, 2011), respectively. After that, the candidate proteins are
64 blasted against a background database to filter out the false-positive records. Then
65 several Perl scripts are used to parse the results to get the putative target protein genes.



66

67

Figure 1. A diagram of the BtToxin_Digger pipeline.

68 3 Results

69 BtToxin_Digger can be easily installed on Linux, macOS, and Windows Subsystem for
70 Linux (WSL) platforms by the conda package manager (Grüning *et al.*, 2018) or docker
71 container. We tested BtToxin_Digger on a laptop with an Intel CPU containing 8 threads
72 of GHz-2.50 and 16 GB memory. It took 14 minutes to process the 1.3-Gbp raw reads
73 to get the results. Moreover, it just takes less than one minute to finish the whole
74 analysis when the other three inputs were provided. BtToxin_Digger can also be used
75 to mine other interesting protein genes with the replacement of the Bt toxin database by
76 other target sequences.

77 We compared BtToxin_Digger with the existing tool BtToxin_scanner (Ye *et al.*, 2012)
 78 and CryProcessor (Shikov *et al.*, 2020). As can be seen from Table 1, BtToxin_Digger
 79 adopts more mining methods, supports more types of input files and toxins, and gets
 80 more friendly output results. Compared with the other two software, it is more suitable
 81 for large-scale toxin gene mining, and at the same time, it can easily implement the
 82 high-throughput analysis.

83 Table 1. Comparison of BtToxin_Digger, BtToxin_scanner and CryProcessor.

Tool	Main methods used	Supported inputs	Supported toxins	Outputs	Flux
BtToxin_Digger	Blast, HMM, SVM	Illumina/Pacbio/Oxford reads, assembled genomes, protein sequences, coding sequences, ORFs	Cry, Cyt, Vip, Other toxins	Toxin list file and sequences file for each input, a matrix file describes all strains vs. all toxins, an integrated file contains sequences and information of all inputs	Unlimited number of inputs with a one-line command
BtToxin_scanner	Blast, HMM, SVM	Assembled genomes, protein sequences, ORFs	Cry	Toxin list file and sequences file for each input	One submit at a time
CryProcessor	HMM	Illumina reads, representing genome assembly graph files, protein sequences	Three-domain Cry	A directory containing multiple files for each input	Unlimited number of inputs with additional shell scripts

84

85 Practice with the sample dataset

86 We also provide the sample dataset to demonstrate the usage of BtToxin_Digger
 87 (Supplementary File 1). To use this tool, users should install it on their computers and

88 have a preliminary understanding of Linux. Users can refer to the protocol (Liu *et al.*,
89 2020) to build their bioinformatics analysis platform and refer to
90 https://github.com/BMBGenomics/BtToxin_Digger#installation to install BtToxin_Digger. We
91 also prepared a webpage ([https://github.com/liaochenlanruo/pgcgap/wiki/Learning-](https://github.com/liaochenlanruo/pgcgap/wiki/Learning-bioinformatics)
92 [bioinformatics](https://github.com/liaochenlanruo/pgcgap/wiki/Learning-bioinformatics)) for users without Linux skills to learn the basic Linux commands.
93 Because the reads file is too large for upload and download, here we only demonstrate
94 the running method of assembled genome, protein sequences, and coding sequences.
95 Users can visit https://github.com/BMBGenomics/BtToxin_Digger#examples for more
96 information.

97 Step 1. Download the Example dataset (Supplementary File 1) and unzip files.

98 Step 2. Open a terminal and enter the directory.

```
99 cd ExampleDataset
```

100 Step 3. Processing assembled genomes

```
101 BtToxin_Digger --SeqPath ./Genome --SequenceType nucl --Scaf_suffix .fas --threads 4
```

102 Step 4. Processing protein sequences

```
103 BtToxin_Digger --SeqPath ./AAs --SequenceType prot --prot_suffix .faa --threads 4
```

104 Step 5. Processing coding sequences

```
105 BtToxin_Digger --SeqPath ./CDSs --SequenceType orfs --orfs_suffix .ffn --threads 4
```

106

107 The running results are stored in Supplementary File 2. *.list: toxin list of each strain;
108 *.gbk: toxin sequences in Genbank format of each strain; Bt_all_genes.table: a matrix
109 describes Strains vs. Toxins; All_Toxins.txt: a table containing all information and
110 sequences of all toxin genes. See Supplementary Table 1 for details.

111

112 **Funding**

113 This work was supported by the National Key R&D Program of China
114 (2017YFD0201201), National Natural Science Foundation of China (31670085,
115 31970003 and 31770003).

116

117 **References**

- 118 Camacho, C., *et al.* (2009) BLAST+: architecture and applications. *BMC Bioinformatics*, 10, 421.
- 119 Chang, C.-C. and Lin, C.-J. (2011) LIBSVM: A library for support vector machines. *ACM Trans. Intell. Syst.*
120 *Technol.*, 2, Article 27.
- 121 Eddy, S.R. (2011) Accelerated Profile HMM Searches. *PLoS Comp Biol*, 7, e1002195.
- 122 G Brookes and Barfoot, P. GM crops: global socio-economic and environmental impacts 1996-2016. UK:
123 PG Economics Ltd; 2018.
- 124 Grüning, B., *et al.* (2018) Bioconda: sustainable and comprehensive software distribution for the life
125 sciences. *Nat Methods*, 15, 475-476.
- 126 Liu, H., *et al.* (2020) Build a bioinformatics analysis platform and apply it to routine analysis of microbial
127 genomics and comparative genomics. *Protocol Exchange*, DOI: 10.21203/rs.2.21224/v2.
- 128 Méric, G., *et al.* (2018) Lineage-specific plasmid acquisition and the evolution of specialized pathogens
129 in *Bacillus thuringiensis* and the *Bacillus cereus* group. *Mol Ecol*, 27, 1524-1540.
- 130 Palma, L., *et al.* (2014) *Bacillus thuringiensis* toxins: an overview of their biocidal activity. *Toxins*, 6, 3296-
131 3325.
- 132 Prado, J.R., *et al.* (2014) Genetically Engineered Crops: From Idea to Product. *Annu Rev Plant Biol*, 65,
133 769-790.
- 134 Ruan, L., *et al.* (2015) Are nematodes a missing link in the confounded ecology of the entomopathogen
135 *Bacillus thuringiensis*? *Trends Microbiol*, 23, 341-346.
- 136 Sanahuja, G., *et al.* (2011) *Bacillus thuringiensis*: a century of research, development and commercial
137 applications. *Plant Biotechnol J*, 9, 283-300.
- 138 Shikov, A., *et al.* (2020) No More Tears: Mining Sequencing Data for Novel Bt Cry Toxins with
139 CryProcessor. *Toxins*, 12, 204.
- 140 Siegwart, M., *et al.* (2015) Resistance to bio-insecticides or how to enhance their sustainability: a review.
141 *Front Plant Sci*, 6, 381.
- 142 Stajich, J.E., *et al.* (2002) The Bioperl toolkit: Perl modules for the life sciences. *Genome Res*, 12, 1611-
143 1618.
- 144 Ye, W., *et al.* (2012) Mining new crystal protein genes from *Bacillus thuringiensis* on the basis of mixed
145 plasmid-enriched genome sequencing and a computational pipeline. *Appl Environ Microbiol*,
146 78, 4795-4801.
- 147 Zheng, J., *et al.* (2017) Comparative Genomics of *Bacillus thuringiensis* Reveals a Path to Specialized
148 Exploitation of Multiple Invertebrate Hosts. *MBio*, 8, e00822-00817.

149