*Research Article*

# Identification and Analysis of Novel Amino-Acid Sequence Repeats in *Bacillus anthracis* str. *Ames* Proteome Using Computational Tools

**G. R. Hemalatha, D. Satyanarayana Rao, and L. Guruprasad**

*School of Chemistry, University of Hyderabad, Hyderabad 500 046, India*

We have identified four repeats and ten domains that are novel in proteins encoded by the *Bacillus anthracis* str. *Ames* proteome using automated in silico methods. A "repeat" corresponds to a region comprising less than 55-amino-acid residues that occur more than once in the protein sequence and sometimes present in tandem. A "domain" corresponds to a conserved region with greater than 55-amino-acid residues and may be present as single or multiple copies in the protein sequence. These correspond to (1) 57-amino-acid-residue PxV domain, (2) 122-amino-acid-residue FxF domain, (3) 111-amino-acid-residue YEFF domain, (4) 109-amino-acid-residue IMxxH domain, (5) 103-amino-acid-residue VxxT domain, (6) 84-amino-acid-residue ExW domain, (7) 104-amino-acid-residue NTGFIG domain, (8) 36-amino-acid-residue NxGK repeat, (9) 95-amino-acid-residue VYV domain, (10) 75-amino-acid-residue KEWE domain, (11) 59-amino-acid-residue AFL domain, (12) 53-amino-acid-residue RIDVK repeat, (13) (a) 41-amino-acid-residue AGQF repeat and (b) 42-amino-acid-residue GSAL repeat. A repeat or domain type is characterized by specific conserved sequence motifs. We discuss the presence of these repeats and domains in proteins from other genomes and their probable secondary structure.

## 1. INTRODUCTION

The anthrax is a disease of herbivores and other mammals including humans, caused by the *Bacillus anthracis* str. *Ames*, a Gram-positive, rod-shaped, nonmotile, spore-forming bacterium [1]. It is an endospore-forming bacterium that causes inhalational anthrax. During the course of disease, endospores are taken up by alveolar macrophages where they germinate in the phagolysosomal compartment. Vegetative cells then escape from the macrophage, eventually infecting blood. Expression of the major plasmid-encoded virulence determinants, tripartite toxin, and a poly-D-glutamic acid capsule is essential for full pathogenicity [2]. Key virulence genes found on plasmids are pXO1 and pXO2 [1]. The 60 MDa plasmid pXO2 carries genes required for the synthesis of an antiphagocytic poly-D-glutamic acid capsule [3]. The 110 MDa plasmid pXO1 [4] is required for the synthesis of the anthrax proteins, edema factor, lethal factor, and protective antigen. These proteins act in binary combinations to produce two anthrax toxins: edema toxin (a protective antigen and edema factor) and lethal toxin (a protective antigen and lethal factor) [5]. The chromosome encodes potential virulence factors that include haemolysins, enterotoxins, phospholipases, proteases, metalloproteases, and iron-acquisition proteins.

The chromosome of *B. anthracis* str. *Ames* contains three homologues of sortase transpeptidase that is responsible for attachment of secreted proteins to peptidoglycan on the cell surface of Gram-positive bacteria [6]. A range of important surface proteins, including enzymes and virulence-related MSCRAMMs (microbial surface components recognizing adhesive matrix molecules) are anchored to the cell wall in Gram-positive bacteria by sortase, a transpeptidase in *Staphylococcus aureus*, that cleaves polypeptides at a conserved LPxTG motif near the carboxyl terminus and covalently links them to penta-glycine crossbridges in peptidoglycan [7, 8]. Nearly 34 candidate surface proteins which have sortase attachment sites and SLH domains were identified. Two putative *B. anthracis* str. *Ames* sortase attached genes have internalin like repeats [9]. The chromosome of

*B. anthracis* str. *Ames* also contains the *csaAB* genes for binding of proteins with S-layer homology (SLH) domains to polysaccharide. The SLH domain is a repetitive modular element that is present in several bacterial cell surface proteins and is involved in noncovalent association with peptidoglycan associated polymers [10]. The SLH domain comprises 55-amino-acid residues [11] and the potential role of most proteins with SLH domains on the surface of *B. anthracis* str. *Ames* is unknown at present [12]. However, these surface proteins may mediate unknown interactions between *B. anthracis* str. *Ames* and its external environment and could be targets for vaccine and drug design. Read et al. [12] reported the complete genome sequence of *B. anthracis* str. *Ames*. It comprises 5 227 293 base pairs and 5508 genes with an overall G+C content of 35.4%. Of these, 2762 are functional genes, 1212 are conserved hypothetical genes, 657 genes are of unknown function, and 877 genes are annotated as hypothetical proteins.

As the complete genome sequence of *B. anthracis* str. *Ames* is available [12], we intended to systematically identify and analyze all the amino-acid sequence repeats in this proteome. In a general context, a "repeat" corresponds to a region comprising less than 55-amino-acid residues that occur more than once, sometimes in tandem along the primary sequence, examples are the YVTN repeats in various cell surface proteins and the WD repeats present in proteins that perform a variety of functions. On the other hand, a "domain" refers to a region of the protein comprising greater than 55-amino-acid residues and does not contain internal sequence repeats. According to the crystallographer definition, a domain represents a region of the protein capable of folding independently as a stable unit. A domain can also exist in multiple copies and there can be several different domains per protein, examples are the SH2, SH3, and PH domains present in signal transduction proteins. The repeats and domains are characterized by conserved sequence motifs that may be identified according to the conservation of individual amino-acid residues at equivalent positions derived from multiple sequence alignments. In the absence of experimental data, the structural information can be obtained from secondary structure or fold prediction studies in silico. Information about the identified domains and repeats is represented in databases such as SMART, INTERPRO and PFAM. SMART (simple modular architecture research tool) allows the identification and annotation of genetically mobile domains and the analysis of domain architectures [13]. INTERPRO is a searchable database that provides information on sequence, function, and annotation. It is an integrated documentation resource for protein families, domains, and sites [14]. PFAM is a large collection of multiple sequence alignments and hidden Markov models covering many common protein domains and families. This can be used to view the domain organization of proteins [15]. We believe that a systematic sequence analysis will provide information on the novel repeats and domains present in *B. anthracis* str. *Ames* proteome that are not identified so far.

The *B. anthracis* str. *Ames* proteome consists of several known repeats and domains. Some of these domains are as follows. (1) BRCT (breast cancer carboxy terminal) domain was first identified as 100-amino-acid tandem repeat at the C-terminus of the tumor suppressor gene product BRCA1, in which the germline mutations lead to nearly 50% familial breast cancer. Most BRCT domains containing proteins participate in DNA damage checkpoint or DNA repair pathways and transcription regulation [16]. The BRCT is an evolutionarily conserved module that exists in a large number of proteins from prokaryotes to eukaryotes. (2) Excalibur (extracellular calcium binding) domain consists of a conserved DxDxDGxxCE motif, which is strikingly similar to the $Ca^{2+}$ binding loop of the calmodulin like EF hand domains, suggesting an evolutionary relationship. (3) Cna_B domain forms a stalk in *Streptococcus aureus* collagen-binding protein that presents the ligand binding domain away from the bacterial cell surface. (4) CBS (cystathionine beta synthase) domain is a small intracellular module with 60-amino-acid residues, mostly found in two or four copies within a protein and occurs in several proteins in all kingdoms of life. Tandem pairs of CBS domains can act as binding domains for adenosine derivatives. In some cases, CBS domains may act as sensors of cellular energy status by being activated by AMP and inhibited by ATP. (5) Par B (par B like nuclease) domain cleaves single stranded DNA, nicks supercoiled plasmid DNA, and exhibits $5'$-$3'$ exonuclease activity. (6) KH (K homology) domain comprises 70-amino-acids residues and is involved in RNA binding. (7) PAS and PAC domains comprising 300 and 45-amino-acid residues, respectively, mediate signal transduction. (8) PASTA domain is an extracellular module comprising 70-amino-acids residues that fold into a globular architecture consisting of 3$\beta$-strands and an $\alpha$-helix which aids in penicillin binding. (9) NEAT (near transporter) domain is a 125-amino-acid residue conserved region consisting mainly $\beta$-strands. The NEAT domain appears to be associated with iron transport in several Gram-positive species, some of them are pathogenic. (10) SLH domain is present in several bacterial cell surface proteins and is involved in noncovalent association with peptidoglycan associated polymers. It comprises 55-amino-acid residues and the predicted secondary structure comprises two $\alpha$-helices flanking a short $\beta$-strand [11].

The repeats present in *B. anthracis* str. *Ames* proteome are as follows. (1) RHS repeats are 21-amino-acids residues long and are involved in carbohydrate binding. (2) TPR (tetratricopeptide) repeats are 34-amino-acids residues long and are involved in protein-protein interactions. (3) EZ_HEAT repeats are 37–47-amino-acid residues long and occur in tandem in a number of cytoplasmic proteins that are involved in intracellular transport processes. Arrays of HEAT repeats consist of 3 to 36 units forming a rod-like helical structure and appear to function as protein-protein interaction surfaces. (4) Ankyrin repeats are about 33-amino-acid residues long and occur in at least four consecutive copies; the core of the repeat appears as a helix-loop-helix structure and is involved in protein-protein interactions. (5) LRR (leucine rich repeats) are 20-amino-acids residues long, each repeat consists of a $\beta$-strand and $\alpha$-helix, that are oriented in an antiparallel manner. The function of LRRs includes signal

transduction, transmembrane receptors, DNA repair, cell adhesion, and extracellualr matrix proteins [17].

Andrade et al. [18] reviewed methods to identify repeats in proteins and the relationship between repeat sequences and their associated functions. Repeats may be identified by manual examination, if the sequence similarity is very high and present in tandem. Repeats are thought to arise due to gene duplication and recombination events. Protein domains may exist either as single or multiple copies and repeats always exist as multiple copies [18, 19]. Programs such as BLASTP [20] are also useful in detecting internal and homologous repeats in a protein database. By using the BLAST program, the presence of repeats in a query protein sequence can be identified if (a) the same region of the query is aligned against two or more distinct regions of a second protein; and (b) different regions of the query are being aligned against the same region of a second protein [18].

Several web-based methods are available for ab initio identification of sequence repeats in proteins. For example, RADAR (rapid automatic detection and alignment of repeats) [21] uses an automatic algorithm, for segmenting a query sequence into repeats; it identifies short composition biased as well as gapped approximate repeats and complex repeat architectures involving many different types of repeats in a query sequence. Rep program [22] uses an iterative algorithm based on score distributions from profile analysis. This procedure allows the identification of homologues at alignment scores lower than the highest optimal alignment score for nonhomologous sequences. The PROSPERO program [23] is ideal for large scale self-comparison of protein sequences. It uses a formula that accurately assesses the significance of protein repeat similarities, allowing for existence of gaps, and also takes into account sequence length and composition. TRUST (tracking repeats using significance and transitivity) program [24] exploits the concept of transitivity of alignments as well as a statistical scheme optimized for the evaluation of repeat significance. Starting from significant local suboptimal alignments, the application of transitivity allows to (1) identify distant repeat homologues for which no alignments were found; (2) gain confidence about consistently well-aligned regions; and (3) recognize and reduce the contribution of nonhomologous repeats. This assessment step will enable to derive a virtually noise-free profile representing a generalized repeat with high fidelity. It has been demonstrated by the authors that TRUST is a useful and reliable tool for mining tandem and nontandem repeats in protein sequence databases, to predict multiple repeat types with varying intervening segments within a single sequence. Once statistically significant repeats are detected, construction of a multiple sequence alignment provides insight into the extent of sequence homology among members of the new protein family and identification of the conserved sequence motifs.

We have implemented TRUST on a personal computer in our laboratory and used it to identify amino-acid sequence repeats in the proteins of *B. anthracis* str. *Ames* proteome. We have identified four repeats and ten domains that are novel in the proteome of *B. anthracis* str. *Ames*. Further analysis

corresponding to searches of the completed and unfinished genome databases identified some of these to be present in other bacterial genomes.

## 2. METHODS

We have downloaded the entire proteome of *B. anthracis* str. *Ames* from the website http://www.ncbi.nlm.nih.gov in the FASTA format. The TRUST program was downloaded from the website and installed on the local Pentium IV computers on the Linux platform. The TRUST server together with the source code is available at http://ibivu.cs.vu.nl/programs/trustwww. The TRUST program was run for all the sequences in this proteome. Based on the size of the TRUST output file, the protein sequences with no internal repeats were discarded automatically; that is, only those protein sequences which comprise repeats were retained. The lengths of repeats and domains currently annotated in the INTERPRO database often comprise greater than 25-amino-acid residues; therefore, in this work, we have considered the repeats with greater than 25-amino-acid residues alone for further analysis. Thus selected proteins were submitted to SMART online (http://smart.embl-heidelberg.de/smart/batch.pl) [13] program in batch mode. Manual inspections of the SMART results identified proteins comprising known repeats or domains and were therefore discarded. Only those repeats that are not identified by SMART database are retained for further analysis.

We have downloaded NCBI NR (release date: April 22, 2005) and UNIPROT (release date: April 23, 2005) databases and installed BLAST-2.2.10 on the local Linux computers (OS: Fedora Core-2, Pentium-IV 3.00 GHz, 1 GB RAM, 80 GB hard disk). Using automatic shell scripts, these protein sequences were then blasted using PSI-BLAST program [25] for three iterations against the NCBI NR database and using BLASTALL program against UNIPROT database. The proteins confirmed to comprise repeats by the BLAST program were retained and were tested for presence in the offline versions of INTERPRO (Database: iprscan_DATA_10.0, Applications: iprscan_V4.1, iprscan_binn4.x_Linux) and PFAM (release date: April 26, 2005) databases. A final check was made using online versions of INTERPRO and PFAM. These series of steps are given in the flowchart as shown in Figure 1.

The repeats which are not present in any of these databases were considered to be novel repeats or domains, depending upon (1) the number of times they occur in the protein sequences, and (2) length of the amino-acid sequence region. The novel repeats and domains thus identified in *B. anthracis* str. *Ames* proteome were subjected to PSI-BLAST analysis in order to identify other proteins from databases that comprise these repeats and domains. Multiple sequence alignment program, ClustalW [26], was used to detect the extent of sequence conservation and the secondary structure prediction was carried out using PHD [27] method.

## 3. RESULTS AND DISCUSSION

From the analysis of *B. anthracis* str. *Ames* proteome using TRUST program, we identified 905 proteins comprising
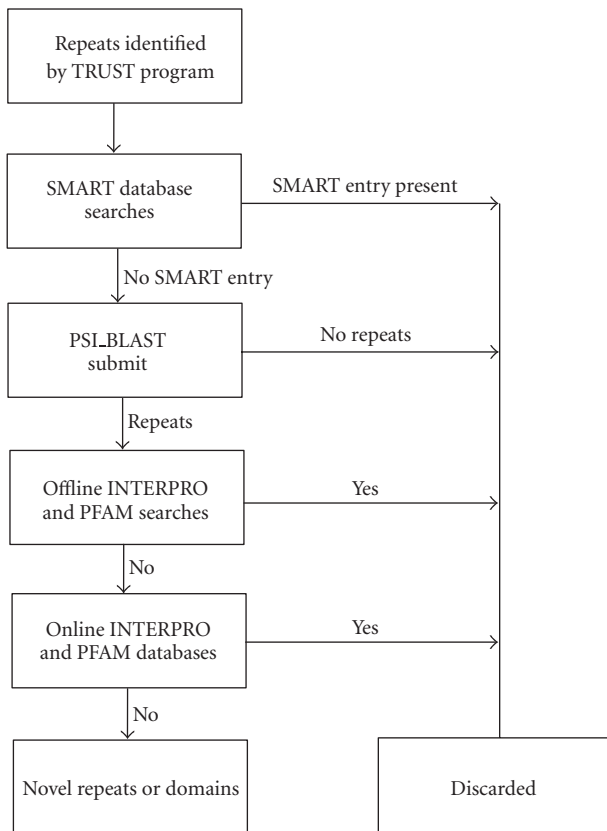
FIGURE 1: Flowchart for systematic analysis of repeats in proteins.

of amino-acid sequence repeats. SMART database analysis identified that 302 entries do not have a SMART description. Based on their absence in the INTERPRO and PFAM databases and the length of repeat sequence (greater than 25-amino-acid residues), we have identified about 120 proteins (data not shown) in the *B. anthracis* str. *Ames* proteome to comprise novel amino-acid sequence repeats. We have added an additional constraint that the repeats identified by TRUST program should also be identified as a repeat by the BLAST program. Subsequent online INTERPRO and PFAM searches confirmed that these domains and repeats have not been reported before. In this work, we have identified four repeats and ten domains, that are not within or part of previously reported repeats and our findings are therefore novel. Further analysis identified some of these in the proteins of other bacterial genomes. The conserved amino-acid residues observed from multiple sequence alignments using the CLUSTALW program were used to describe sequence motifs characteristic of these novel repeats and domains. Often, more than one sequence motif is associated with repeats or domains and the amino-acid sequence patterns characteristic of these repeats are represented according to the PROSITE description [28]. Ponting et al. [29], have earlier used a similar approach to identify novel domains and repeats in *Drosophila melanogaster*.

In this work, we identified four repeats and ten domains that have not been reported before in the *B. anthracis* str.

*Ames* proteome. The repeats and domains described in 1 to 6 and 9 are also present in some bacterial organisms, 7, 8, 10 and 11 are *Bacillus*-specific, 12 and 13 are *Bacillus anthracis* str. *Ames* specific. Lists of the proteins containing these novel repeats and domains are shown in Tables 1a to 1k. These tables indicate the protein identifiers (Gene or Swall_ID), the number of amino-acid residues in the protein, a description of the protein, and other well-characterized repeats and domains present in the protein. Some sequences representing these repeats or domains share lower than 15% pairwise sequence identity. However, these sequences retain the conserved motifs and the positions of secondary structure elements in the multiple sequence alignment. For all the proteins, the amino-acid sequence corresponding to each representative repeat are shown in the multiple sequence alignments (see Figures from 2 to 14).[1] Conservation of the position of secondary structural elements is indicated from the multiple sequence alignment. The schematic figures used to represent these repeats and domains are shown in Figures 15 to 27. These figures (drawn to an approximate scale) reflect the relative proximity and location of individual repeats and domains along the primary sequence. We discuss each of these novel repeats and domains below.

### 3.1. 57-amino-acid-residue PxV domain

The 251-amino-acid-residue protein corresponding to the GENE_ID BA2292 and described as hypothetical protein comprises of a 57-amino-acid-residue region as two copies. Further BLAST searches using sequence corresponding to the region (65–121) as a query identified 24 proteins that are described as hypothetical (see Table 1(a)). This region occurs as four copies in proteins from *Shewanella amazonensis, and Haloarcula marismortui*, as two copies in proteins from *B. anthracis, B. cereus, B. halodurans, B. thuringiensis, B. thuringiensis serovar, Thermus thermopilus, Chloroflexus aurantiacus, Chloroflexus aggregans Exiguobacterium sp., Bacillus weihenstephanensis, Roseiflexus castenholzii, Clostridium novyi, Herpetosiphon aurantiacus*, and as single copy in *Anabaena variabilis*; we therefore describe this region as a

---

[1] The multiple sequence alignments corresponding to representative repeats and domains from various proteins along with their GENE or SWall identifiers. (a) PxV domain, (b) FxF domain, (c) YEFF domain, (d) IMxxH domain, (e) VxxT domain, (f) ExW domain, (g) NTGFIG domain, (h) NxGK repeat (i) VYV domain, (j) KEWE domain, (k) AFL domain, (l) RIDVK repeat, (m)(a) AGQF repeat and (b) GSAL repeat. The numbers given in brackets indicate the start and end of amino-acid-residue positions corresponding to either the repeat or domain. The 80% consensus is labeled according to the alignment to the alignment generated at the website http://www.bork.embl-heidelberg.de/Alignment/consensus.html: alcohol (o, ST); aliphatic (I, ILV); any (·, ACDEFGHIKLMNPQRSTVWY); aromatic (a, FHWY); charged (c, DEHKR); hydrophobic (h, ACFGHIKLM-RTVWY); negative (−, DE); polar (p, CDEHKNQRST); positive (+, HKR); small (s, ACDGNPSTV); tiny (u, AGS); turn-like (t, ACDEGHKN-QRST). A capital letter indicates 80% conservation of corresponding amino-acid residue. The secondary structure prediction indicated at the top was derived using the PHD program. Residues predicted with greater than 82% accuracy to form β-sheets are represented by "E" and α-helices are represented by "H."

TABLE 1: The proteins are represented by their corresponding Gene_ID along with the number of amino-acid residues indicated in brackets in the first column. The organism and corresponding phylogeny are indicated in the second column: (A) represents Archaea and (B) represents Bacteria. The third column contains the description of the proteins containing the repeats or the domains identified elsewhere, including those identified in the present work and the total number of such repeats or domains. The fourth column represents exclusively the total number of novel repeats or domains identified in this work.

(a) List of proteins containing the 57-amino-acid-residue PxV domain.

| Gene ID (number of residues) | Organism | Description | Number of PxV domains |
|---|---|---|---|
| BA2292 (251) | *Bacillus anthracis* str. *Ames* (B) | Hypothetical protein | 2 |
| BAS2138 (249) | *Bacillus anthracis* Sterne (B) | Hypothetical protein | 2 |
| BT9727_2076 (249) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | Hypothetical protein | 2 |
| BCZK2072 (249) | *Bacillus cereus* E33L (B) | Hypothetical protein | 2 |
| BCE2326 (249) | *Bacillus cereus* ATCC 10987 (B) | Hypothetical protein | 2 |
| BC2244 (249) | *Bacillus cereus* ATCC 14579 (B) | Hypothetical protein | 2 |
| BH1282 (222) | *Bacillus halodurans* C-125 (B) | BH1282 protein | 2 |
| BCE_G9241_2259 (249) | *Bacillus cereus* G9241 (B) | Hypothetical conserved protein | 2 |
| RBTH_03198 (251) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Hypothetical protein | 2 |
| TT_P0044 (221) | *Thermus thermopilus* HB27 (B) | Hypothetical conserved protein | 2 |
| TTHB089 (221) | *Thermus thermophilus* HB8 (B) | Hypothetical protein | 2 |
| Chlo02001630 (262) | *Chloroflexus aurantiacus* J-10-fl (B) | Hypothetical protein | 2 |
| ExigDRAFT_0608 (264) | *Exiguobacterium sibiricum* 255-15 (B) | Hypothetical protein | 2 |
| SamaDRAFT_3539 (469) | *Shewanella amazonensis* SB2B (B) | Hypothetical protein | 4 |
| rrnAC0576 (488) | *Haloarcula marismortui* ATCC 43049 (A) | Unknown | 4 |
| Ava_3757 (292) | *Anabaena variabilis* ATCC 29413 (B) | Hypothetical protein | 1 |
| BcerKBAB4DRAFT_2942 (249) | *Bacillus weihenstephanensis* KBAB4 (B) | Conserved hypothetical protein | 2 |
| B14911_22687 (254) | *Bacillus* sp. NRRL B-14911 (B) | Hypothetical protein | 2 |
| Bcer98DRAFT_2673 (249) | *Bacillus cereus* subsp. cytotoxis NVH (B) | Conserved hypothetical protein | 2 |
| RcasDRAFT_0590 (259) | *Roseiflexus castenholzii* DSM 13941 (B) | Surface protein from Gram-positive cocci, anchor region | 2 |
| RoseRSDRAFT_1732 (259) | *Roseiflexus* sp. RS-1 (B) | Surface protein from Gram-positive cocci, anchor region | 2 |
| NT01CX_1619 (210) | *Clostridium novyi* NT (B) | Conserved hypothetical protein | 2 |
| HaurDRAFT_2803 (196) | *Herpetosiphon aurantiacus* ATCC 23779 (B) | Conserved hypothetical protein | 2 |
| CaggDRAFT_2922 (261) | *Chloroflexus aggregans* DSM 9485 (B) | Conserved hypothetical protein | 2 |

TABLE 1: Continued.

(b) List of proteins containing the 122-amino-acid-residue FxF domain.

| Gene ID (number of residues) | Organism | Description | Number of FxF domains |
|---|---|---|---|
| BA0881 (293) | *Bacillus anthracis* str. *Ames* (B) | Conserved domain protein | 2 |
| BCZK0785 (293) | *Bacillus cereus* E33L (B) | Hypothetical protein | 2 |
| BT9727_0783 (295) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | Hypothetical protein | 2 |
| BCE_G9241_0886 (293) | *Bacillus cereus* G9241 (B) | Conserved protein, putative | 2 |
| GK3171 (297) | *Geobacillus kaustophilus* HTA426 (B) | Hypothetical conserved protein | 2 |
| CTC00525 (279) | *Clostridium tetani* E88 (B) | Hypothetical protein | 2 |
| Bcer98DRAFT_3031 (293) | *Bacillus cereus* subsp. cytotoxis NVH (B) | Conserved hypothetical protein | 2 |
| B14911_04439 (305) | *Bacillus* sp. NRRL B-14911 (B) | Hypothetical protein | 2 |
| DredDRAFT_0533 (262) | *Desulfotomaculum reducens* MI-1 (B) | Hypothetical protein | 2 |
| NT01CX_1557 (276) | *Clostridium novyi* NT (B) | Conserved protein, putative | 2 |

(c) List of proteins containing the 111-amino-acid-residue YEFF domain.

| Gene ID (number of residues) | Organism | Description and other known domains | Number of YEFF domains |
|---|---|---|---|
| BA3695 (510) | *Bacillus anthracis* str. *Ames* (B) | S-layer protein, putative, SLH-domain (3) | 2 |
| BCZK3337 (492) | *Bacillus cereus* E33L (B) | S-layer protein, SLH-domain (3) | 2 |
| BT9727_3386 (510) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | S-layer protein, SLH-domain (3) | 2 |
| Bant_01004347 (510) | *Bacillus anthracis* str. A2012 (B) | Hypothetical protein, SLH-domain (3) | 2 |
| BCE_G9241_3590 (492) | *Bacillus cereus* G9241 (B) | Lipoprotein, putative SLH-domain (3) | 2 |
| BA5326 (321) | *Bacillus anthracis* str. *Ames* (B) | Lipoprotein, putative | 2 |
| BT9727_4791 (321) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | Hypothetical protein | 2 |
| BC5098 (321) | *Bacillus cereus* ATCC 14579 (B) | Hypothetical protein | 2 |
| BCZK4809 (321) | *Bacillus cereus* E33L (B) | Hypothetical protein | 2 |
| RBTH_06214 (321) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Hypothetical protein | 2 |
| EF0374 (325) | *Enterococcus faecalis* V583 (B) | Lipoprotein, putative | 2 |
| EF0375 (321) | *Enterococcus faecalis* V583 (B) | Hypothetical protein | 2 |
| EF0376 (347) | *Enterococcus faecalis* V583 (B) | Hypothetical protein | 2 |

(d) List of proteins containing the 109-amino-acid-residue IMxxH domain.

| Gene ID (number of residues) | Organism | Description | Number of IMxxH domains |
|---|---|---|---|
| BA1021 (266) | *Bacillus anthracis* str. *Ames* (B) | Hypothetical protein | 2 |
| BAS0955 (283) | *Bacillus anthracis* Sterne (B) | Hypothetical protein | 2 |
| BCZK0933 (283) | *Bacillus cereus* E33L (B) | Hypothetical protein | 2 |
| BT9727_0941 (283) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | Hypothetical protein | 2 |
| BC1029 (283) | *Bacillus cereus* ATCC 14579 (B) | Hypothetical protein | 2 |
| RBTH_03050 (283) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Hypothetical protein | 2 |
| CAC3450 (307) | *Clostridium acetobutylicum* ATCC 824 (B) | Hypothetical protein | 2 |
| CPE0158 (303) | *Clostridium perfringens* str. 13 (B) | Hypothetical protein | 2 |
| CTC02189 (314) | *Clostridium tetani* E88 (B) | Conserved protein | 2 |
| CtheDRAFT_1311 (307) | *Clostridium thermocellum* ATCC 27405 (B) | Conserved hypothetical protein | 2 |
| DhafDRAFT_0725 (321) | *Desulfitobacterium hafniense* DCB-2 (B) | Conserved hypothetical protein | 2 |

TABLE 1: Continued.

(d) Continued.

| Gene ID (number of residues) | Organism | Description | Number of IMxxH domains |
|---|---|---|---|
| BCE_G9241_1042 (283) | *Bacillus cereus* G9241 (B) | Conserved protein | 2 |
| CbeiDRAFT_3331 (312) | *Clostridium beijerincki* NCIMB 8052 (B) | Conserved hypothetical protein | 2 |
| CphyDRAFT_3436 (305) | *Clostridium phytofermentans* ISDg (B) | Conserved hypothetical protein | 2 |
| ClosDRAFT_1658 (308) | *Clostridium* sp. OhILAs (B) | Conserved hypothetical protein | 2 |
| CdifQ_02001573 (254) | *Clostridium difficile* QCD-32g58 (B) | Hypothetical protein | 2 |
| BcerKBAB4DRAFT_3543 (283) | *Bacillus weihenstephanensis* KBAB4 (B) | Hypothetical protein | 2 |
| AmetDRAFT_1908 (272) | *Alkaliphilus metalliredigenes* QYMF (B) | Conserved hypothetical protein | 2 |
| CD1511 (304) | *Clostridium difficile* 630 (B) | Conserved hypothetical protein | 2 |
| CPF_0149 (303) | *Clostridium perfringens* ATCC 13124 (B) | Hypothetical protein | 2 |
| BcerKBAB4DRAFT_0307 (171) | *Bacillus weihenstephanensis* KBAB4 (B) | Conserved hypothetical protein | 1 |
| Bcer98DRAFT_1038 (303) | *Bacillus cereus* subsp. cytotoxis NVH 391-98 (B) | Conserved hypothetical protein | 2 |

(e) List of proteins containing the 103-amino-acid-residue VxxT domain.

| Gene ID (number of residues) | Organism | Description | Number of VxxT domains |
|---|---|---|---|
| BA4716 (349) | *Bacillus anthracis* str. *Ames* (B) | Germination protein gerM | 2 |
| gerM BT9727_4219 (349) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | Germination protein | 2 |
| germ BCZK4235 (349) | *Bacillus cereus* E33L (B) | Germination protein | 2 |
| BCE4587 (349) | *Bacillus cereus* ATCC 10987 (B) | Germination protein gerM | 2 |
| BC4495 (349) | *Bacillus cereus* ATCC 14579 (B) | Germination protein germ | 2 |
| BSU28380 (366) | *Bacillus subtilis* subsp. subtilis str. 168 (B) | Germination protein gerM | 2 |
| BL00314 (369) | *Bacillus licheniformis* ATCC 14580 (B) | Spore germination protein GerM | 2 |
| BH3070 (365) | *Bacillus halodurans* C-125 (B) | Germination (Cortex hydrolysis) and sporulation | 2 |
| RBTH_05210 (349) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Germination protein germ | 2 |
| gerM (210) | *Bacillus subtilis* (B) | gerM | 1 |
| ABC2653 (377) | *Bacillus clausii* KSM-K16 (B) | Germination protein GerM | 2 |
| GK2667 (357) | *Geobacillus kaustophilus* HTA426 (B) | Germination (Cortex hydrolysis) and sporulation | 2 |
| OB2107 (352) | *Oceanobacillus iheyensis* HTE831 (B) | Germination (Cortex hydrolysis) and sporulation | 2 |
| SwolDRAFT_2302 (195) | *Syntrophomonas wolfei* str. Goettingen (B) | Hypothetical protein | 1 |
| MothDRAFT_0979 (200) | *Moorella thermoacetica* ATCC 39073 (B) | Similar to Spore germination protein | 1 |
| CtheDRAFT_0840 (299) | *Clostridium thermocellum* ATCC 27405 (B) | Hypothetical protein | 1 |
| gerM ABF83609 (349) | *Bacillus thuringiensis* serovar kurstaki (B) | Spore germination protein | 2 |
| Bcer98DRAFT_3179 (348) | *Bacillus cereus* subsp. cytotoxis NVH 391-98 (B) | Germination protein GerM | 2 |
| BcerKBAB4DRAFT_4089 (349) | *Bacillus weihenstephanensis* KBAB4 (B) | Germination protein gerM | 2 |
| B14911_06091 (361) | *Bacillus* sp. NRRL B-14911 (B) | Spore germination protein | 2 |
| GAA01614 (295) | *Pelotomaculum thermopropionicum* SI (B) | Unnamed protein product | 1 |
| AmetDRAFT_1640 (332) | *Alkaliphilus metalliredigenes* QYMF (B) | Hypothetical protein | 2 |
| Moth_0516 (200) | *Moorella thermoacetica* ATCC 39073 (B) | Spore germination protein-like | 1 |

TABLE 1: Continued.

(f) List of proteins containing the 84-amino-acid-residue ExW domain.

| Gene ID (number of residues) | Organism | Description | Number of ExW domains |
|---|---|---|---|
| BA4310 (246) | *Bacillus anthracis* str. *Ames* (B) | Hypothetical protein | 2 |
| BT9727_3829 (246) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | Hypothetical protein | 2 |
| BCE4157 (246) | *Bacillus cereus* ATCC 10987 (B) | Hypothetical protein | 2 |
| BCZK3845 (246) | *Bacillus cereus* E33L (B) | Hypothetical protein | 2 |
| BC4088 (248) | *Bacillus cereus* ATCC 14579 (B) | IG hypothetical 17224 | 2 |
| GK0969 (226) | *Geobacillus kaustophilus* HTA426 (B) | Hypothetical conserved protein | 2 |
| BSU30660 (145) | *Bacillus* subtilis subsp. str. 168 (B) | Hypothetical protein ytkA (PSPA8) | 1 |
| BL05305 (147) | *Bacillus licheniformis* ATCC 14580 (B) | Conserved protein YtkA | 1 |
| BH0983 (157) | *Bacillus halodurans* C-125 (B) | BH0983 protein | 1 |
| Bant_01004966 (252) | *Bacillus anthracis* str. A2012 (B) | Protein chain release factor A | 2 |
| RBTH_02670 (248) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Hypothetical protein | 2 |
| BCE_G9241_4093 (246) | *Bacillus cereus* G9241 (B) | IG hypothetical protein | 2 |
| OB2488 (166) | *Oceanobacillus ihenyensis* HTE831 (B) | Hypothetical conserved protein | 1 |
| ABC0230 (158) | *Bacillus clausii* KSM-K16 (B) | Unknown conserved protein | 1 |
| BH0678 (246) | *Bacillus halodurans* C-125 (B) | BH0678 protein | 2 |
| ABC4088 (142) | *Bacillus clausii* KSM-K16 (B) | Hypothetical protein | 1 |
| ExigDRAFT_1796 (161) | *Exiguobacterium sibiricum* 255-15 (B) | Hypothetical protein | 1 |
| OB3282 (155) | *Oceanobacillus ihenyensis* HTE831 (B) | Hypothetical conserved protein | 1 |
| BcerKBAB4DRAFT_2040 (241) | *Bacillus weihenstephanensis* KBAB4 (B) | Conserved hypothetical protein | 2 |
| B14911_09907 (144) | *Bacillus* sp. NRRL B-14911 (B) | Hypothetical protein | 1 |
| B14911_05359 (273) | *Bacillus* sp. NRRL B-14911 (B) | Hypothetical protein | 2 |
| BAA83944 (267) | *Bacillus halodurans* (B) | Unnamed protein product | 2 |
| BH1853 (158) | *Bacillus halodurans* C-125 (B) | Hypothetical protein | 1 |
| Bcer98DRAFT_3614 (177) | *Bacillus cereus* subsp. cytotoxis NVH 391-98 (B) | IG hypothetical protein | 2 |
| ExigDRAFT_0574 (253) | *Exiguobacterium sibiricum* 255-15 (B) | Hypothetical protein | 2 |

(g) List of proteins containing the 104-amino-acid-residue NTGFIG domain.

| Gene ID (number of residues) | Organism | Description | Number of NTGFIG domains |
|---|---|---|---|
| BA2665 (232) | *Bacillus anthracis* str. *Ames* (B) | Hypothetical protein | 2 tandem |
| BT9727_2444 (232) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | Hypothetical protein | 2 tandem |
| BCZK2413 (232) | *Bacillus cereus* E33L (B) | Group-specific protein | 2 tandem |
| BCE2700 (234) | *Bacillus cereus* ATCC 10987 (B) | Hypothetical protein | 2 tandem |
| BC2674 (234) | *Bacillus cereus* ATCC 14579 (B) | Hypothetical protein | 2 tandem |
| Bant_01003317 (236) | *Bacillus anthracis* str. A2012 (B) | Hypothetical protein | 2 tandem |
| BCE_G9241_CNI_0263 (234) | *Bacillus cereus* G9241 (B) | Conserved hypothetical protein | 2 tandem |
| BcerKBAB4DRAFT_0535 (232) | *Bacillus weihenstephanensis* KBAB4(B) | Conserved hypothetical protein | 2 tandem |
| Bcer98DRAFT_0128 (234) | *Bacillus cereus* subsp. cytotoxis NVH 391-98 (B) | Conserved hypothetical protein | 2 tandem |

TABLE 1: Continued.

(h) List of proteins containing the 36-amino-acid-residue NxGK repeat.

| Gene ID (number of residues) | Organism | Description and other known domains | Numbre of NxGK repeats |
|---|---|---|---|
| BA3686 (193) | *Bacillus anthracis* str. *Ames* (B) | Hypothetical protein, SAP domain (1) | 2 |
| BT9727_3378 (193) | *Bacillus thuringiensis* serovar konkukian str. 97-27 (B) | Hypothetical protein, SAP domain (1) | 2 |
| BCZK3328 (193) | *Bacillus cereus* E33L (B) | Hypothetical protein, SAP domain (1) | 2 |
| BC3626 (193) | *Bacillus cereus* ATCC 14579 (B) | Hypothetical protein, SAP domain (1) | 2 |
| BCE3645 (193) | *Bacillus cereus* ATCC 10987 (B) | Hypothetical protein, SAP domain (1) | 2 |
| RBTH_03615 (193) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Hypothetical cytosolic protein, SAP domain (1) | 2 |
| BCE_G9241_3579 (193) | *Bacillus cereus* G9241 (B) | Hypothetical cytosolic protein, SAP domain (1) | 2 |
| BcerKBAB4DRAFT_0944 (193) | *Bacillus weihenstephanensis* KBAB4 (B) | Conserved hypothetical protein, SAP domain (1) | 2 |
| B14911_25780 (189) | *Bacillus* sp. NRRL B-14911 (B) | Hypothetical protein, SAP domain (1) | 2 |

(i) List of proteins containing the 95-amino-acid-residue VYV domain.

| Gene ID (number of residues) | Organism | Description | Number of VYV domains |
|---|---|---|---|
| BA1701 (225) | *Bacillus anthracis* str. *Ames* (B) | Hypothetical protein | 2 tandem |
| BAS1577 (227) | *Bacillus anthracis* str. Sterne (B) | Hypothetical protein | 2 tandem |
| RBTH_03882 (1004) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Hypothetical exported protein | 10 tandem |
| DSY3134 (1674) | *Desulfitobacterium hafniense* Y51 (B) | Hypothetical protein | 2 tandem |

(j) List of proteins containing the 75-amino-acid-residue KEWE domain.

| Gene ID (number of residues) | Organism | Description | Number of KEWE domains |
|---|---|---|---|
| BA3147 (262) | *Bacillus anthracis* str. *Ames* (B) | Hypothetical protein | 3 tandem |
| BAS2924 (344) | *Bacillus anthracis* str. Sterne (B) | Hypothetical protein | 4 tandem |
| RBTH_06405 (331) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Hypothetical protein | 4 tandem |
| pE33L466_0092 (328) | *Bacillus cereus* E33L (B) | Hypothetical protein | 4 tandem |
| Bant_01003795 (178) | *Bacillus anthracis* str. A2012 (B) | Hypothetical protein | 2 tandem |
| pBMB165 (247) | *Bacillus thuringiensis* serovar tenebrionis (B) | Hypothetical protein | 3 tandem |

(k) List of proteins containing the 59-amino-acid-residue AFL domain.

| Gene ID (number of residues) | Organism | Description | Number of AFL domains |
|---|---|---|---|
| BA3065 (290) | *Bacillus anthracis* str. *Ames* (B) | Hypothetical protein. | 2 |
| BAS2851 (297) | *Bacillus anthracis* str. Sterne (B) | Hypothetical protein | 2 |
| Bant_01003715 (293) | *Bacillus anthracis* str. A2012 (B) | Hypothetical protein | 2 |
| RBTH_02124 (145) | *Bacillus thuringiensis* serovar israelensis ATCC 35646 (B) | Hypothetical protein | 1 |
| BcerKBAB4DRAFT_1832 (291) | *Bacillus weihenstephanensis* KBAB4 (B) | Conserved hypothetical protein | 2 |

domain. The length of proteins varied between 196 to 488-amino-acid residues. The multiple sequence alignment corresponding to this domain is associated with PxV sequence motif where x is any amino-acid residue and is shown in Figure 2. The pairwise identities between sequences corresponding to PxV domain varied between 15–96%. The secondary structure corresponding to PxV domain is predicted to comprise four $\beta$-strands as shown in Figure 2. The representative domain architecture corresponding to proteins comprising the PxV domain is shown in Figure 15.

```
Secondary structure              EEEEE            EE                                EEEEE        EE

RcasDRAFT_0590_1(32-89)      VRVIHAS-PDAPAVDVIVNGNR--ALTNVPFFAASAYLDLPAGSYDIQVVPAGAT-S-PVVID    58
RoseRSDRAFT_1732_1(32-89)    VRVVHAS-PDAPAVDVIVNGNK--ALTNVPFFAASAYLDLPAGSYDIQVVPAGAT-S-PVVID   58
Chlo02001630_1(32-90)        VRVIHAS-PDAPAVDVFVNGNA--VLTNVGFFAASPYLDLPAGTYRVQVAPTGAG-AGSAVID   59
CaggDRAFT_2922_1(31-89)      VRVIHAS-PDAPAVDVFVNGNA--VLTNVGFFAASPYLDLPAGTYRVQVAPTGAG-AGSAVID   59
HaurDRAFT_2803_1(4-62)       VRVMHAS-PDAPAVDIFVDGKA--VLTSVPFFALSGQLALPDGTYTIDIAPAGAG-VAASVFE   59
B14911_22687_1(67-124)       VRVVHAS-PDAPNVDIYVNGNR--ILKDFPYKDVSGYLSLPAGKYQIDIYPAGDM-V-STVLS   58
HaurDRAFT_2803_2(105-162)    VRVIHGS-PDAPAVDIKIAGTQN-VVVKGADFGDAATLEVPAGTYSFDISPAGSS-T--VLFT   58
rrnAC0576_1(67-124)          VRVAHMS-PNAPNVDVYLEGDA--VLEDVPFGAVSQYLDVPAGERSVEITAAGD--PDTSVFS   58
rrnAC0576_2(284-341)         VRVAHMS-PNAPNVDVYVDGSA--VLEDVPFGAVSDYLEVPAGARTVEITAAGD--PDTSVFE   58
BH1282_1(30-89)              VRVLHAS-PDAPPVDVYIDGKK--QMEGVPFKQTSSYFNVPAGDHMITIFAAGDDPAETPVIE   60
ExigDRAFT_0608_1(29-86)      VRVIHAS-PDAPAVDIAVDGKK--AVSGAEFKAVTDYLTLPAGEHKVEVFAAGT--TKDPVLS   58
RBTH_03198_2(161-218)        IRFAHFS-PDTPVVNVDLKDGDH-LFENVLFKQITDFLQVSPGTADIEISLANNK---NVLLT   58
BC2244_2(159-216)            IRFAHFS-PDTPVVNVDLKDGDH-LFENVLFKQITDFLQVSPGTADIEISLADNK---NVLLT   58
BcerKBAB4DRAFT_2942_2(159-216) IRFAHFS-PDTPVVNVNLKDGDH-LFENVLFKQITDFLQVSPGTADIEVSLADTK---KVLLT  58
BCE_G9241_2259_2(159-217)    IRFAHFS-PDTPVVNVSLKGGDH-LFENVLFKQITDFLEVSPGTADIEVSLADNQ---NVLLT   58
BCE_2326_2(159-216)          IRFAHFS-PDTPVVNVSLKGGDH-LFENVLFKQITDFLEVSPGTADIEVSLADHQ---SVLLT   58
BA2292_2(161-218)            IRFAHFS-PDTPVVNVSLKDGDH-LFENVLFKQITDFLEVSPGTADIEVSLADNQ---SVLLT   58
BAS2138_2(159-216)           IRFAHFS-PDTPVVNVSLKDGDH-LFENVLFKQITDFLEVSPGTADIEVSLADNQ---SVLLT   58
BT7727_2076_2(159-216)       IRFAHFS-PDTPVVNVSLKDGDH-LFENVLFKQITDFLEVSPGTADIEVSLADNQ---SVLLT   58
BCZK2072_2(159-216)          IRFAHFS-PDTPVINVSLKDGDH-LFENVLFKQITDFLEVSPGTADIEVSLADNQ---SILLT   58
Bcer98DRAFT_2673_2(159-216)  IRFAHFS-PDTSVVNVSLKNGDH-LFENVLFKQVTDYLQVSPGTADIEISLADTK---KNLVT   58
NT01CX_1619_2(113-170)       VKFVHLS-PGTPNVDITLPNGTI-LFKDVEFEEGTDYIPLKVGTYTIEAKPTGSD---KTVLT   58
B14911_22687_2(164-221)      ARFIHLS-PDAPAVDIAVKKGDV-IFPNISFRQATQYLGLTPMTVDLEVRVAGSS---NTVLS   58
RcasDRAFT_0590_2(131-188)    VRVIHFS-PDAPAVDIKVAGGPT-LISNLAFPNASNYLPVDAGSYDLQVTPAGGT---AVVLD   58
RoseRSDRAFT_1732_2(131-188)  VRVIHFS-PDAPAVDIKVAGGPT-LISNLAFPNASNYLPVDAGSYDLQVTPAGGT---AVVLD   58
Chlo02001630_2(131-188)      VRVYHFS-PDAPAVDVKLANGTT-LISNLAFPNASDYLEVPAGTYDLQVTPAGGS---AVVIN   58
CaggDRAFT_2922_2(130-187)    VRVYHFS-PDAPAVDVKLANGTT-LISNLAFPDASDYLEVPAGTYDLQVTPAGGD---AVVIN   58
ExigDRAFT_0608_2(126-183)    VRVAHFA-PDAPAVDVAPKGGDP-LFSDLEFSKVSDYGTLDAGTYDLEVRPAGAT---DVVKA   58
TTHB089_2(124-181)           IRVVHAS-PDAPAVDVAVKGGPV-LFAGLPFPRASAYASVPAGTYDLEVRAAGTA---TVALD   58
TTP0044_2(124-181)           IRVVHAS-PDAPAVDVAVKGGPV-LLAGLPFPRASAYASVPAGTYDLEVRAAGTA---TVALD   58
BH1282_2(130-187)            LRAVHLS-PDTPAVQLHLSAANV-DMPSLSFENASRYIDLPAGAYDLDIRMIETD--DVATE   58
RBTH_03198_1(65-121)         IRIFHAD-PNIPAVDILVNGQKV--IKNISFKQFSPYLSLVQGKYRIDIVPVGNET---PIFS   57
BC2244_1(63-119)             IRIFHAD-PNIPAVDILVNGQKV--IKNISFKQFSPYLSLVQGKYRIDIVPVGNET---PIFS   57
BCE_G9241_2259_1(63-119)     IRFFHSA-SNTPAVDILVNGQKV--IKNISFKQFSPYLTLVQGKYRIDIVPVGNET---PIFS   57
BAS2138_1(63-119)            IRFFHSA-SNTPAVDILVNGQKV--IKNISFKQFSPYLTLVQGKYRIDIVPVGNET---PIFS   57
BCE_2326_1(63-119)           IRFFHSA-SNTPAVDILVNGQKV--IKNISFKQFSPYLTLVQGKYRIDIVPVGNET---PIFS   57
BA2292_1(65-121)             IRFFHSA-SNTPAVDILVNGQKV--IKNISFKQFSPYLTLVQGKYRIDIVPVGNET---PIFS   57
BCZK2072_1(63-119)           IRFFHSA-SNTPAVDILVNGQKV--IKNISFKQFSPYLTLVQGKYRIDIVPVGNET---PIFS   57
BT7727_2076_1(63-119)        IRFFHSA-SNTPAVDILVNGQKV--IKNISFKQFSPYLTLVQGKYRIDIVPVGNET---PIFS   57
BcerKBAB4DRAFT_2942_1(63-119) MRIFHTA-PHTPAVDIIINGQKV--IKNISFKQFSPYLSLVQGKYRIDIVPVGNET---PIFS  57
Bcer98DRAFT_2673_1(63-119)   MRIFHAS-PHTPAVDILVNGQKV--IKNITFQQFSPYFSLMQGQYRLDIVPLDNET---PIFS   57
SamaDRAFT_3539(264-321)      IRVAHSA-ADVPQVDILANGTKVDALSGAAFGQASGYLNLAPGEYQVDTVLTSDNS---VVGI   59
Ava_3757(63-120)             LRVINAAVPTASPVDVIVNGQRV--LENVNFRQASRYVNVTPGNIQVLFVTSGTNS---TIAS   58
TTHB089_1(24-81)             VRVAHLS-PDAPAVDVLVNGQRA--ITGLAFKEVTPYIPLPAAKVRVQVVPAGQDAP--VVID   58
TTP0044_1(24-81)             VRVAHLS-PDAPAVDVLVNGQRA--ITGLAFKEVTPYIPLPAAKVRVQVVPAGQDAP--VVID   58
NT01CX_1619_1(15-72)         MRLLNAS-PNAPAVDVYFNGQLI--TSNLAYKEFTEYMSTSPGLYNVKVFPHGKLSS--PIID   58

consensus/80%                lRhhHhu.PssPsVsl.lpstt...hpsl.F.phosalpls.Gphplpl..ssst....slhs
```

Figure 2: BA2292 is homologous to protein GBAA2292 from *Bacillus anthracis* str. "*Ames* Ancestor." BAS2138 is homologous to proteins BT9727_2076 from *Bacillus thuringiensis* serovar konkukian str. 97-27 and Bant_01002917 from *Bacillus anthracis* str. A2012.

### 3.2. 122-amino-acid-residue FxF domain

The 293-amino-acid-residue protein corresponding to the GENE_ID BA0881 and described as conserved domain protein comprises a 122-amino-acid-residue region as two copies. Further BLAST searches using sequence corresponding to the region (55–176) as a query identified 10 proteins (see Table 1(b)). The proteins comprising this region are described as either conserved or hypothetical proteins. This region occurs as two copies in the proteins of *B. anthracis, B. cereus, B. thuringiensis, Geobacillus kaustophilus, Clostridium tetani, Clostridium novyi,* and *Desulfotomaculum reducens* genomes. The length of proteins varied between 262 to 305-amino-acid residues. The multiple sequence alignment corresponding to this domain is associated with characteristic sequence motif FxF (Figure 3) and we refer to this as the FxF domain. The pairwise sequence identities corresponding to this domain varies between 18–97%. The secondary structure corresponding to FxF domain is predicted to comprise one α-helix and five β-strands, and the representative domain architecture of proteins comprising this domain is shown in Figure 16.

### 3.3. 111-amino-acid-residue YEFF domain

The 510-amino-acid-residue protein corresponding to the GENE_ID BA3695 and described as a S-layer protein comprises a 111-amino-acid-residue region that is present as two copies. Further BLAST searches, using sequence corresponding to the region (247–357) as a query, identified 13 proteins (see Table 1(c)), that are described as S-layer proteins, hypothetical, or lipoproteins and correspond to the *B. anthracis* str. *Ames* and A2012, *B. cereus, B. thuringiensis, B. thuringiensis* serovar israelensis, and *Enterococcus faecalis* genomes. The length of proteins varied between 321 to 510-amino-acid residues. Five proteins corresponding to the GENE_ID BA3695 and Bant_01004347 of *B. anthracis,* BCE_G9241_3590, and BCZK3337 of *B. cereus* and BT9727_3386 of *B. thuringiensis* comprise three copies of SLH domain, indicating a cell surface role for these proteins.

```
Secondary structure        HHHHHHH            EEEE                EEEEE               EEEE          EEEE

BA0881_1(55-176)           I Y Q F L H K E L P R L E E Y Q I S L S G I E I E K R D N G - Y D V A V F I R S T V P K P I S F E E V T L I L L N K E K K L C A R K T
BCZK0785_1(55-176)         I Y Q F L H K E L P R L E E Y Q I S L S G I E I E E R D N G - Y D V A V F I R S T V P K P I S F E E V T L I L L N K E K K L C A R K T
BCE_G9241_0886_1(55-176)   I Y Q F L H K E L P R L E E Y Q I S L S G I E I E K R D S G - Y D V A V F I R S T V P K P I S F E E V T L I I L L N K E K K L C A R K T
Bcer98DRAFT_3031_1(55-176) I Y Q F L H K E L P R L Q E N Q I S L S G I E I E K R E G S - Y A V A A F I R S S I S K P I S F E E V T L L L L N K E D E L C A R K T
BT9727_0783_1(58-179)      I Y Q F L H K S L P T L Q E N Q I S L A G I E S K K H E N A - Y Y I T T F I R S S V K H P I Q F E T L T L S L L N K N G E T C A R Q T
GK3171_1(46-167)           V Y R F Y H E Q L P P L Q P N Q I S I S G V K L V E Y N D G - F V A V A I L R N T L P K P V R F E R I R L L L L D E D G T A I A R K E
B14911_04439_1(59-182)     V L R F L N N E L P P L L P N Q I S L A G I E L Q Q D G G S - L T V A A F V R S S L S K A V E F K K T H L L L V G P D E E I L A R K E
BA0881_2(185-293)          A L R N F V D N L T P P N D G E I N F L G L Q A A R K E N G D L H T T L L I R N G C K D N I Q L E Q L P L H I E D A T G A V V V K G A
BCZK0785_2(185-293)        A L R N F V D N L T P P N D G E I N F L G L Q A A R K E N G D L H T T L L I R N G C K D N I Q L E Q L P L H I E D A T G A V V V K G A
BCE_G9241_0886_2(185-293)  A L R N F V D N L T P P N D G E I N F L G L Q A A R K E N G D L H T T L L I R N G C K E N I Q L E Q L P L H I E D A T G A V V V K G A
Bcer98DRAFT_3031_2(185-293) A L R N F V E S L T P P Q N G E L N F L G L Q A A Q K E N G D L H A T I L I R N G C K R N I Q L K Q L P L H I E D A S G E I V V K G A
BT9727_0783_2(188-295)     K L Q E I I A N L D P P E E D E I N F R G L N A V V E E N G D L N A T I L I R N G Y N K N I T L E Q L P L H I S D R S E S T V A E R I
B14911_04439_2(191-305)    K L K Q M V E Q M D P P K I G E I N F M G I Q A K V A D N E D L Q V T L L I R N G N D Q N V M L Q Q L P L Q V E D A T S E V I A K G G
GK3171_2(176-297)          Q L Q A L V D S V P P P A P G E V N F M G I E A K Q L P S G E L G V T L L I R N G S D K H I H F E Q I P L E V R D Y A G D I V A R G L
NT01CX_1557_2(164-276)     Q Y E K F L K E L P L L R E G Q V T M N A Y D V Y T N E D D G I A V E L V I R N G R H N G V D I K R I P L S I Y D K D K K L V A S G T
DredDRAFT_0533_2(156-262)  Q F T T F L K K L P S V Q E G S I N I D T Y S I E K N N D G S L T V A I V L R H R L A K P T V L S R F Q F G I V D T N K S I V A R A A
CTC00525_2(170-279)        V F K E F L E S L P K L E R G Q G S I S V F T I T Q Y E N G D L L M T L L V R N A T D E A V T M T K M P I T L K T Q K G E T I L S G V
CTC00525_1(36-159)         L E E E L R E V I P K V E E G K I N I A G I Y A F D Q G D K - V E V K A Y L A N G L S Q K I N F E D V P I Y I I N S K E E K L A Y Q V
NT01CX_1557_1(31-154)      C L E E E L E A L P A I K E G E L D V N - V D F F F D L G D R Y E A S I F I R N G L S T G V N L E K I P F I V L D K D E K E V G R K I
DredDRAFT_0533_1(25-147)   L M Q E E I N N L P Q I T D G T V A I D S I Y T V N W E D K - I E I G F Y L R N V T S H K I C F T Q T P L K I L N P K G E V L A S V T

consensus/80%              h h p . h h c p L s . . p p s p l s h . u l p h . p t p s s . h t s s h h l R s s h t c s l p h c p l s L h l . s t p t p h h s c t h
```

```
Secondary structure        EEEEE

BA0881_1(55-176)           F N L S A L G D I P S N V N M P F I F T F E Q E T I - T D A A L S Q T D W E L A F E L E S K - - H T L D L D P S W E A    122
BCZK0785_1(55-176)         F N L S A L G D I P A N V N M P F I F T F E Q E T I - T D A A L S Q T D W E L A F E L E S K - - H A L D L D P S W E A    122
BCE_G9241_0886_1(55-176)   F N L S A L G D I P A N V N M P F I F T F E Q E T I - T D A D L S Q T D W E L A F E L E S K - - H V L D L D P S W E A    122
Bcer98DRAFT_3031_1(55-176) F N L S D I G D I P A N V N M P W V F T F D E E T I - T D A E L S Q T D W Q L A F E L E G E - - H R L D L D P T W E T    122
BT9727_0783_1(58-179)      F D L S H L E G I P S N V N M P W T F V F E E N S I - T E A T L S N E D W Q L V F E L Q G K - - H S L D L D P I W Q E    122
GK3171_1(46-167)           F D M S P F G E L P P M T A R P W R F L F A A E D K - L V D Q L P A D G W K I A F E L T P R - - H R L D L E E S W E Q    122
B14911_04439_1(59-182)     F D L T E I G E I P A K S S R P W N F T F N S S D L - L T D S I P A E G W K L A F E I R N N E E H R L D L D E A W E N    124
BA0881_2(185-293)          F T L P N L E I K A N - T T K P W S F V F P A S S I - L K E D M D L S S W K A L V P Q D - - - - - - - - - - - - - - -    109
BCZK0785_2(185-293)        F T L P N L E I K A N - T T K P W S F V F P A S S I - L K E D M D L S S W K A L V P Q D - - - - - - - - - - - - - - -    109
BCE_G9241_0886_2(185-293)  F T L P N L E I K A N - T T K P W S F V F P A S S I - L K E D M E L S S W K A L V P Q D - - - - - - - - - - - - - - -    109
Bcer98DRAFT_3031_2(185-293) F T L P N L E I K A N - S T K P W S F I F P V S F V - L K K E M D L S T W K A I V P Q D - - - - - - - - - - - - - - -    109
BT9727_0783_2(188-295)     F V L K D F Q I K A N - S T K P W T F T F P A D S V - S K E P I D L S K W K A F I P Q - - - - - - - - - - - - - - - -    108
B14911_04439_2(191-305)    F Q L D K F E L K A N - T S K P W T F I F P K S L L - L K D N P D L S S W K A Y P L Q Q Q V Q T E I - - - - - - - - -    115
GK3171_2(176-297)          F P C H - L E V K A H - T S K P W T F L F P P E L L - H K A E P D W T S W K V T I P S S P A Q S E K Q E T P S S D E -    122
NT01CX_1557_2(164-276)     F Y L E D A S L N P I - S A K V Y L F T F S K D E L - L R E D Y N L K N W T I Q F L L N S N V N - - - - - - - - - - -    113
DredDRAFT_0533_2(156-262)  F V I E Q Y I L E P G - M F L L R S F K F T P E T I - V N S D A D I N Q C S I A F L - - - - - - - - - - - - - - - - -    107
CTC00525_2(170-279)        F D I E N F T V N P Y - K A R V L S L I F K K E V V N I E E D F D L S T C K I I F E R E - - - - - - - - - - - - - - -    110
CTC00525_1(36-159)         F D L S E E G D I P S G K A I P V K L N F N K Q N I - L V D K I P Q D D W K V V F G G N D V K G V R Y V N I E L E S I    124
NT01CX_1557_1(31-154)      F N L R E V G E I P A R S V R P W K I Y F E K D E L - N V E G I N L K D L K I V F D S R I K A A G V V N V Q Y E N L P    124
DredDRAFT_0533_1(25-147)   I N L S D M G D I P A Y S V R P W R F Y L G K E D L - - T L D N S L K D L K I A F N S R N I P P Y M L V I E D R L P E    123

consensus/80%              F s L p t h t . h s s . s s h P a . F h F . t p p l . h p t p h s . p s W c h h h . . p . . . . . . . . . . . . . . .
```

FIGURE 3: BA0881 is homologous to proteins GBAA0881 *Bacillus anthracis* str. "*Ames* Ancestor," BAS0837 from *Bacillus anthracis* str. Sterne and Bant_01001534 from *Bacillus anthracis* str. A2012.

This domain is characterized by conserved sequence motifs; YEFF, RGD, FTY, GKD, and FVEH. We refer to this 111-amino-acid region as the YEFF domain. The pairwise sequence identities corresponding to the YEFF domain varied between 36–96%. The consensus secondary structure predicted for this domain suggests mainly $\beta$-strands and the conserved sequence motifs, that is, YEFF and FTY are associated with $\beta$-strands; see Figure 4. The representative domain architecture of proteins comprising this domain is shown in Figure 17. It is intriguing that each domain comprises RGD sequence motif which is found in the proteins of extracellular matrix. Many viruses enter their host cells via the RGD motif—integrin interaction and synthetic peptides containing this RGD motif are active modulators of cell adhesion [30]. The RGD motif was originally identified as the sequence within fibronectin that mediates cell attachment. This motif has now been found in numerous other proteins and supports cell adhesion. The integrins, a family of cell surface proteins, act as receptors for cell adhesion molecules. A subset of the integrins recognizes the RGD motif within their ligands, the binding of which mediates both cell substratum and cell-cell interactions [31]. The presence of RGD motif and SLH domain implies that the YEFF domain compris-ing proteins is also present on the cell surface and mediates protein-protein interactions.

### 3.4. 109-amino-acid-residue IMxxH domain

The 266-amino-acid-residue protein corresponding to the GENE_ID BA1021 and described as hypothetical protein comprises a 109-amino-acid-residue region as two copies. Further BLAST searches using sequence corresponding to the region (4–112) as a query identified 22 proteins (see Table 1(d)) that are described as either conserved or hypothetical proteins. This domain region occurs as two copies in all the proteins of *B. anthracis*, *B. cereus*, *B. thuringiensis, Bacillus weihenstephanensis C. acetobutylicum, C. perfringens, C. tetani, C. thermocellum, Desulfitobacterium hafniense, Clostridium phytofermentans*, and *Alkaliphilus metalliredigenes*, and as single domain in the 171-amino-acid-residue protein BcerKBAB4DRAFT_0307. The length of proteins varied between 171 to 321-amino-acid residues. The multiple sequence alignment corresponding to this domain identified the characteristic sequence motifs; IMxxH, REA, and we refer to this as the IMxxH domain. The IMxxH sequence motif occurs at the N-terminal region of the domain. The

```
Secondary structure              E E E E               E E E E     E E E E            E E E
EF0374(62-172)       I L S S - - T D W Q G T K V Y D K N N N N L T A E N A N F I G L A K Y D G E T G F Y E F F D K E T G E T R G D E G T F F V T D - - - D G E
EF0375(58-168)       I L S G - - T D W Q G T R V Y D A A G N D L T A E N A N F I G L A K Y D G E T G F Y E F F D K N T G E T R G D E G T F F V T G - - - D G T
EF0376(59-172)       G L S E - - K D W A G T R V Y D R N G N D L T D E N Q N L L H A I K F D A T T S F Y E F F D K E T G E S T G D E G T F F M T A G I T D V S
BA5326(58-168)       I L S D - - T N W Q G T R V L D K D K N D V T K E N A N F I G L A K Y D A K S G R Y E F F D A K T G A S R G D K G T F F I T N - - - D G K
BCZK4809(58-168)     I L S D - - T N W Q G T R V L D K D K N D V T K E N A N F I G L A K Y D A K S G R Y E F F D A K T G A S R G D K G T F F I T N - - - D G K
BT9727_4791(58-168)  I L S D - - T N W Q G T R V Y D K D K N D V T K E N A N F I G L A K Y D A K S G R Y E F F D A K T G A S R G D K G T F F I T N - - - D G K
BC5098(58-168)       I L S E - - T N W Q G T R V Y D K D K N D V T K E N A N F I G L A K Y D A K S G R Y E F F D A K T G A S R G D K G T F F V T N - - - D G K
RBTH_06214(58-168)   I L S K - - T N W Q G T R V Y D K D K N D L T K E N A N F I G L A K Y D A K S G R Y E F F D A K T G A S R G D K G T F F V T N - - - D G K
BA3695(247-357)      I L G E - - T N W Q G T K V Y D K D H N D V T K E N Q N F I G L A K Y D A K T A R Y E F F N A S T G E S R N D S G T F F I T N - - - D G K
Bant_01004347(247-357) I L G E - - T N W Q G T K V Y D K D H N D V T K E N Q N F I G L A K Y D A K T A R Y E F F N A S T G E S R N D S G T F F I T N - - - D G K
BT9727_3386(247-357) I L G E - - T N W Q G T K V Y D K D H N D V T K E N Q N F I G L A K Y D A K M A R Y E F F N A S T G E S R N D S G T F F I T N - - - D G K
BCZK3337(229-339)    I L G E - - T N W Q G T K V Y D K D H N D V T K E N Q N F I G L A K Y D A K T A R Y E F F N A S T G E S R N D S G T F F I T N - - - D G K
BCE_G9241_3590(229-339) I L G E - - T N W Q G T K V Y D K D H N D V T K E N Q N F I G L A K Y D A K T A R Y E F F N A K T G E S R N D S G T F F I T N - - - D G K
EF0376(223-336)      F D G T P Q L L W N G T K V V D K D G N D V T S A N Q N F I S L A K F D Q D S S K Y E F F N L Q T G E T R G D Y G Y F K V G N - - - Q N K
EF0375(199-310)      I L G T - - T L W N G T K V V D K N G N D V T A A N Q N F I S L A K F D P N T S K Y E F F N L Q T G E T R G D F G Y F Q V V D - - - N N K
EF0374(203-314)      I L G A - - T L W N G T K V L D E D G N D V T E A N K M F I S L A K F D N K T S K Y E F F D L E T G K T R G D F G Y F Q V I D - - - N N K
BA3695(388-499)      I L S S - - T L W N G T V V L D E Q G N N V T K Y N S N L I S L A K Y D K N T N K Y E F F N V N T G E S R G D Y G F F D V V H - - - D N K
BT9727_3386(388-499) I L S S - - T L W N G T V V L D E Q G N N V T K Y N S N L I S L A K Y D K N T N K Y E F F N V N T G E S R G D Y G F F D V V H - - - D N K
Bant_01004347(388-499) I L S S - - T L W N G T V V L D E Q G N N V T K Y N S N L I S L A K Y D E N T N K Y E F F N V N T G E S R G D Y G F F D V V H - - - D N K
BCZK3337(370-481)    I L S S - - T L W N G T V V L D E Q G N N V T K Y N S N L I S L A K Y D K N T N K Y E F F N V N T G E S R G D Y G F F D V V H - - - G N K
BCE_G9241_3590(370-481) I L S S - - T L W N G T V V L D D Q G N D V T K Y N S N L I S L A K Y D K N T N K Y E F F N V N T G E S R G D Y G F F D V V H - - - G N K
BA5326(199-310)      I L G G - - T L W H G T K V L D E A G N D V T Q F N S N F I S L A K F D D K S N K Y E F F N S E T G Q S R G D Y G Y F D V L H - - - E N K
BCZK4809(199-310)    I L G G - - T L W H G T K V L D E A G N D V T Q F N S N F I S L A K F D D K F N K Y E F F N S E T G Q S R G D Y G Y F D V L H - - - E N K
BT9727_4791(199-310) I L G G - - T L W H G T K V L D E T G N D V T Q F N S N F I S L A K F D D K S N K Y E F F N S E T G Q S R G D Y G Y F D V L H - - - E N K
BC5098(199-310)      I L G G - - T L W H G T K V L D E A G N D V T Q F N S N F I S L A K F D D K S N K Y E F F N S E T G Q S R G D Y G Y F D V V H - - - E N K
RBTH_06214(199-310)  I L G G - - T L W H G T K I L D E A G N D V T Q F N S N F I S L A K F D D K S N K Y E F F N S E T G Q S R G D Y G Y F D V V H - - - E N K

consensus/80%        I L u t . . T . W p G T + V h D c s t N D l T p . N t N h I u L A K a D t p o s + Y E F F s h p T G p S R G D . G h F . l s p . . . - s K
```

```
Secondary structure   E E E E E        E E E E        E E E              E E E E E
EF0374(62-172)       K R I L I S D T Q N - Y Q A V V D L T E V T K D K F T Y K R M G K D K D G K D V E V F V E H I P   111
EF0375(58-168)       K R I L I S R T Q N - Y Q A V V D L T E V S K D K F T Y K R L G K D K L G N D V E V Y V E H I P   111
EF0376(59-172)       R L V I I S E T K N - Y Q G V Y P L R T L Y Q D T F T Y R Q M G K D K N G N D I E V F V E N K A   114
BA5326(58-168)       K R I L I S E S M K - Y Q A V V D M T K L N K N V F T Y K R M G K D A N G N D V E V F V E H V P   111
BCZK4809(58-168)     K R I L I S E S M K - Y Q A V V D M T K L N K N V F T Y K R M G K D A N G N D V E V F V E H V P   111
BT9727_4791(58-168)  K R I L I S E S M K - Y Q A V V D M T K L N K N V I F T Y K R M G K D A N G N D V E V F V E H V P   111
BC5098(58-168)       K R I L I S E S M K - Y Q A V I D M T K L N K N V F T Y K R M G K D A N G K D V E V F V E H V P   111
RBTH_06214(58-168)   K R I L I S E S M K - Y Q A V V D M T K L N K N V F T Y K R M G K D A N G K D V E V F V E H V P   111
BA3695(247-357)      K R V L I S E T Q N - Y Q A V V E L T Q L D K E K F T Y K R M G K D A K R N D V E V F V E H I P   111
Bant_01004347(247-357) K R V L I S E T Q N - Y Q A V V E L T Q L D K E K F T Y K R M G K D A K R N D V E V F V E H I P   111
BT9727_3386(247-357) K R V L I S E T Q N - Y Q A V V E L T Q L D K E K F T Y K R M G K D A K G N D V E V F V E H I P   111
BCZK3337(229-339)    K R V L I S E T Q N - Y Q A V V E L T Q L D K E K F T Y K R M G K D V K G N D V E V F V E H I P   111
BCE_G9241_3590(229-339) K R V L I S E T Q N - Y Q A V V E L T Q L D K E K F T Y K R M G K D V K G N D V E V F V E H I P   111
EF0376(223-336)      F R A H V S I G T N R Y G A V L E L T E L N D N R F T Y T R M G K D N E G N D I Q V Y V E H E P   114
EF0375(199-310)      I R A H V S I G T N R Y G A A L E L T E L N D R F T Y T R M G K D N A G N D I Q V F V E H E P   112
EF0374(203-314)      I R A H V S I G D N K Y G A A L E L T E L N D K R F T Y T R M G K D N N G K E I K V F V E H E P   112
BA3695(388-499)      I R A H V S L G N N K Y G A V L E L T E L N K E K F T Y T R M G K D A N G K D I K I F V E H E P   112
BT9727_3386(388-499) I R A H V S L G N N K Y G A V L E L T E L N K E K F T Y T R M G K D A N G K D I K I F V E H E P   112
Bant_01004347(388-499) I R A H V S L G N N K Y G A V L E L T E L N K E K F T Y T R M G K D A N G K D I K I F V E H E P   112
BCZK3337(370-481)    I R A H V S L G N N K Y G A V L E L T E L N K A K F T Y T R M G K D A N G K D I K I F V E H E P   112
BCE_G9241_3590(370-481) I R A H A S L G N N K Y G A V L E L T E L N K E K F T Y T R I G K D A N G K D I K I F V E H E P   112
BA5326(199-310)      I R A H V S I G N N K Y G A A L E L T E L N K N K F T Y K R T G K D Q A G N D I T I F V E H E P   112
BCZK4809(199-310)    I R A H V S I G N N K Y G A A L E L T E L N K N K F T Y K R T G K D Q A G N D I T I F V E H E P   112
BT9727_4791(199-310) I R A H V S I G N N K Y G A A L E L T E L N K N K F T Y K R T G K D Q A G N D I T I F V E H E P   112
BC5098(199-310)      I R A H V S I G N N K Y G A A L E L T E L N K N K F T Y K R T G K D Q A G N D I T I F V E H E P   112
RBTH_06214(199-310)  I R A H V S I G N N K Y G A A L E L T E L N K N K F T Y K R T G K D Q A G K D I T I F V E H E P   112

consensus/80%        h R h h l S . s p N . Y t A s l - L T p L s K p p F T Y p R h G K D t t G p D l p l F V E H . P
```

Figure 4: BA3695 is homologous to proteins GBAA3695 from *Bacillus anthracis* str. "*Ames* Ancestor" and BAS342 from *Bacillus anthracis* str. *Sterne*. BA5326 is homologous to proteins GBAA5326 from *Bacillus anthracis* str. "*Ames* Ancestor," BAS4948 from *Bacillus anthracis* str. Sterne and Bant_01000199 from *Bacillus anthracis* str. A2012.

pairwise sequence identities corresponding to the IMxxH domain varies between 5–98%. The secondary structure corresponding to IMxxH domain is predicted to comprise four α-helices as shown in Figure 5. The representative domain architecture corresponding to proteins comprising this domain is shown in Figure 18.

### 3.5. 103-amino-acid-residue VxxT domain

The 349-amino-acid-residue protein corresponding to the GENE_ID BA4716 and described as germination protein comprises a 103-amino-acid-residue region as two copies. Further BLAST searches using sequence corresponding to the region (67–169) as query identified 23 proteins (see

Table 1(e)). The proteins comprising this domain are described as germination proteins as the *Bacillus anthracis* is an endospore-forming bacterium. This domain region occurs twice in proteins of *B. anthracis* str. *Ames, B. cereus, B. clausii, B. thuringiensis*, *B. thuringiensis serovar israelensis*, *Alkaliphilus metalliredigene*, and *Bacillus weihenstephanensis* genomes and only once in the proteins of *Syntrophomonas wolfei str. Goettingen, Moorella thermoacetica, Clostridium thermocellum, B. subtilis*, and *Pelotomaculum thermopropionicum* genomes. The length of proteins varied between 195 to 377-amino-acid residues. The multiple sequence alignment corresponding to this domain identified VxxT as sequence motif. This sequence motif occurs in the N-terminal region of each protein and the pairwise sequence identity

```
Secondary structure            HHHHHHHHHHHHHHHHHHH          HHHHHHHHHHHHHHHHHHHHHH

BCE_G9241_1042_1(21-129)       ERSLNEIRFWSRIMKEHSLFLRLGFRCEDTQLIEEANQFYRLFEHIEQIAHSYTNETDPEQ-----IKRF
BCZK0933_1(21-129)             ERSLNEIRFWSRIMKEHSLFLRLGFRCEDTQLIEEANQFYRLFEHIEQIAHSYTNETDPEQ-----IKRF
BT9727_0941_1(21-129)          ERSLNEIRFWSRIMKEHSLFLRLGFRCEDTQLIEEANQFYRLFEHIEQIAHSYTNETDPEQ-----IKRF
BA1021_1(4-112)                ERSLNEIRFWSRIMKEHSLFLRLGFRCEDTQLIEEANQFYRLFEHIEQIAHSYTNETDPEQ-----IKRF
BAS0955_1(21-129)              ERSLNEIRFWSRIMKEHSLFLRLGFRCEDTQLIEEANQFYRLFEHIEQIAHSYTNETDPEQ-----IKRF
RBTH_03050_1(21-129)           ERSLNEIRFWSRIMKEHSFFLRLGFRCEDTQLIEEANQFYRLFEHIEQIAHSYTNETDPEQ-----IKRF
BC1029_1(21-129)               ERSLNEIRFWSRIMKEHSFFLRLGFRCEDTQLIEEANQFYRLFEHIEQIAHSYTNETDPEQ-----IKRF
BcerKBAB4DRAFT_3543_1(21-129)  ERSLNEIRFWSRIMKEHSFFLRLGFRCEDTQLIEEANQFYRLFEHIEQIAYSYTNETDPGQ-----IKRF
Bcer98DRAFT_1038_1(42-147)     EKSLTENRFWLRIMKEHALFLGEGFNRKDTNLIQQVDQFFHLFDRHLQKAFSIP--QTVQA-----VRQL
CTC02189(189-294)              RYAYEQETTWNRIMAEHAKFIRGLLDPTEDALIDTANNFGKEFDELTR---EAKRAMYKTM----PISKV
CbeiDRAFT_3331(190-295)        REAYEQEAFWNRIMAEHSKFIRGLLDPTEDELINTANNFGHQFDILTR---EARAAMNKSI----PISKV
ClosDRAFT_1658(189-294)        KEIYEQELFWNRIMAEHSKFIRGLLDPTEDELIHIANDFAKEFDALTA---AVEEAIEKCL----PIDKI
CtheDRAFT_1311(189-294)        KEAYELQFFWNRQMAEHAKFIRGLLDPTENDLINQANDFGNEFDQLTA---EAKAAMDATS----PMAKV
CdifQ_02001573(138-241)        KNAKEIELFWDHIMMEHALFMRGLLDPSEGELINTSNDFAIKFNELIE---KTN--EMTDS----NIKNI
CD1511(189-291)                KNAKEIELFWDHIMMEHALFMRGLLDPSEGELINTSNDFAIKFNELIE---KTN--EMTDS----NIKNI
CPE0158_2(188-291)             VNISKTEAFWNEIMMEHSLFIRGLLDPSEYELINTAHEFAFEFNELIQ--QLN-NVTNV----TIDNV
CPF_0149(188-291)              VNISKTEAFWNEIMMEHSLFIRGLLDPSEYELINTAHEFAFEFNELIQ--QLN-NVTNV----TIDNV
CphyDRAFT_3436(189-292)        EDLKDDELFWNQIMMEHALFIRGLLDPTENDLIMQADDFASVYADLLD---EAS--TMTER----TMGDL
DhafDRAFT_0725_2(197-302)      CHMVEMQMFWDHIMEHAEVISHLLDPKEKAMITRADHFAQAYEQLLN--QLGNGTVPDQ-----SFRRI
BCZK0933_2(149-260)            DAIIKENVFFLRIMADHAKFIGHLLDPSERKLVDTARNFSNDFDALMYQAIDLESMKPQSQ-TVPLLDQF
BT9727_0941_2(149-260)         DAIIKENVFFLRIMADHAKFIGHLLDPSERKLVDTARNFSNDFDALMYQAIDLESMKPQSQ-TVPLLDQF
BA1021_2(132-243)              DAIIKENVFFLRIMADHAKFIGHLLDPSERKLVDTARNFSNDFDALMYQAIDLESMKPQSQ-TVPLLDQF
BAS0955_2(149-260)             DAIIKENVFFLRIMADHAKFIGHLLDPSERKLVDTARNFSNDFDALMYQAIDLESMKPQSQ-TVPLLDQF
BCE_G9241_1042_2(149-260)      DAIIKENVFFLRIMADHAKFIGHLLDPSERKLVDTARNFSNDFDALMYQAIDLESMKPQSQ-TVPLLDQF
BC1029_2(149-260)              DAIIKENVFFLRIMADHAKFIGHLLDPSERKLVDTARNFSNDFDELMYQAIDLESMKPQSQ-TAPLLDQF
RBTH_03050_2(149-260)          DAIIKENVFFLRIMADHAKFIGHLLDPSERKLVDTARNFSNDFDELMYQAIDLESMKPQSQ-TVPLLDQF
BcerKBAB4DRAFT_3543_2(149-260) DAIIKENVFFLRIMADHAKFIGHLLDPSERKLVDTARNFSNDFDELMYQAIDLESMKPQSQ-TVPLLDQF
BcerKBAB4DRAFT_0307(35-147)    DAIISENVFLRIMMEHSRFIGSLLDQSERNLVHTALKFGDDDFEILLNQARDVESMLYQKEPTYPIIGKM
Bcer98DRAFT_1038_2(167-279)    DAIISENVFWLRIMMEHSRFIASLLDQSERNLVHTALKFGDDFEVLLSQARDVESMLYQKQPTYPIIGKM
CAC3450_1(190-295)             QGIIRQEIFWNDIMEDHAEFIRGYLDPSQTSLFNTANNFVRRFDDIEN---ATESLTNNPS----NLNNI
CPE0158_1(9-119)               TSSLELHLFFMRVMKEHAIFLEAGLGPKNSKLAKELDKCKGNLEKLLFDVVKLSKGRVRQSIVD-SGEVF
DhafDRAFT_0725_1(12-122)       RESLELHLFWARIIKEHLIFLESGFMCKDADWMQEADALKCSFEEILHEANCLADGKVGIEVMK-SGELF
CAC3450_2(9-121)               RLSLELNLFFLRIVKEHNVIAGASLPPKYAPTLMEILAVNKKLDMLLSKTVALSKGNISREAMN-SSTLI
AmetDRAFT_1908_1(11-115)       NVALFEHQFWLQVLGDHARFILNALSPEEREEIQRAQYFIHIFDQLLE----ESRKSPRGS----ALSKL
AmetDRAFT_1908_2(133-245)      TQPIHYHMVWLLDAAGHSAGIMGDLDMVEKELIRKSGKFTQRFEEFYIKAVEIAGYTRTTLDQFPAFTRF

consensus/80%                  ct.hp..hFa.+IMt-HuhFlthhhcsp-ppLlppAppF.p.F-tl.....phpt.p..pp.....lpph


Secondary structure            HHHHHHH   HHHHHHHHH

BCE_G9241_1042_1(21-129)       NAEVQQAATNIWGFKRKILGLILTCKLPGQNNFPLLVDHTSREA  109
BCZK0933_1(21-129)             NAEVQQAATNIWGFKRKILGLILTCKLPGQNNFPLLVDHTSREA  109
BT9727_0941_1(21-129)          NAEVQQAATNIWGFKRKILGLILTCKLPGQNNFPLLVDHTSREA  109
BA1021_1(4-112)                NAEVQQAATNIWGFKRKILGLILTCKLPGQNNFPLLVDHTSREA  109
BAS0955_1(21-129)              NAEVQQAATNIWGFKRKILGLILTCKLPGQNNFPLLVDHTSREA  109
RBTH_03050_1(21-129)           NAEVQQAATNIWGFKRKILGLILTCKLPGQNNFPLLVDHTSREA  109
BC1029_1(21-129)               NAEVQQAATNIWGFKRKILGLILTCKLPGQNNFPLLVDHTSREA  109
BcerKBAB4DRAFT_3543_1(21-129)  NSEVQQAATNIWGFKRKILGLILTCKLPGQNNFPLLVDHTSREA  109
Bcer98DRAFT_1038_1(42-147)     NEESIQLVYAFRNYKRNLLILINCKVSGFN-FPLLVDHIAREA   106
CTC02189(189-294)              TNRSLRATRRIRNFKKQGTEGILDCKIRSII-IPLLADHTLREA  106
CbeiDRAFT_3331(190-295)        TDESLEATKSIRNFKAQGTQGLVECKIKSII-IPLLGDHTLREA  106
ClosDRAFT_1658(189-294)        TDKSLEATKEVRNFNTQGTEGLLDCKIRSII-IPLLGDHVLRES  106
CtheDRAFT_1311(189-294)        TDESLKATEDFRNFKAQGTQAILECKVKSII-IPLLGDHVLREA  106
CdifQ_02001573(138-241)        TEETLNETVEFKDFKEAGASGIEQCKIKSII-LPLLADHVLREA  104
CD1511(189-291)                TEETLNETVEFKDFKEAGASGIEQCKIKSII-LPLLADHVLREA  104
CPE0158_2(188-291)             THEILKETTRLRDFKEEGTKGIMNCNIKSLI-LPLLSDHVLREA  104
CPF_0149(188-291)              THETLKETTRLRDFKEEGTKGIMNCNIKSLI-LPLLSDHVLREA  104
CphyDRAFT_3436(189-292)        TCRTLEETIKYRDFKLAGTKGINDCEIRSII-LPLLADHVLREA  104
DhafDRAFT_0725_2(197-302)      TSETIRVTGEFKDFKAAGTDAILCCQLRSLI-LPLLADHVLREA  106
BCZK0933_2(149-260)            LDQNRVSVASLRDFKKTARDLIEQCKIKSII-HPLLADHVFREA  112
BT9727_0941_2(149-260)         LDQNRVSVASLRDFKKTARDLIEQCKIKSII-HPLLADHVFREA  112
BA1021_2(132-243)              LDQNRVSVASLRDFKKTARDLIEQCKIKSII-HPLLADHVFREA  112
BAS0955_2(149-260)             LDQNRVSVASLRDFKKTARDLIEQCKIKSII-HPLLADHVFREA  112
BCE_G9241_1042_2(149-260)      LDQNRVSVASLRDFKKTARDLIEQCKIKSII-HPLLADHVFREA  112
BC1029_2(149-260)              LDQNRVSVASLRDFKKTARDLIEQCKIKSII-HPLLADHVFREA  112
RBTH_03050_2(149-260)          LDQNRVSVTSLRDFKKTARDLIEQCKIKSII-HPLLADHVFREA  112
BcerKBAB4DRAFT_3543_2(149-260) LDQNRVSVTSLRDFKKTARDLIEQCKIKSII-HPLLADHVFREA  112
BcerKBAB4DRAFT_0307(35-147)    NKDSENATVELRNFKKAGLELIQTCQIRSVI-NPLLADHVTREA  113
Bcer98DRAFT_1038_2(167-279)    NKDSENATVELRNFKKAGLELIQTCQIRNVI-NPLLADHVTREA  113
CAC3450_1(190-295)             TRNIYSLVTEFRNFKSTATKGLLACKIKAIM-APLLADHVTREA  106
CPE0158_1(9-119)               TDYTLETEKKTEHYTGININSKITTMEKDLMC--APKKGIDSKV  111
DhafDRAFT_0725_1(12-122)       TNKTLKAEQKTQELTCIPINSQLTVETMSLHP-YMGVGMGMVP   111
CAC3450_2(9-121)               TPLTLPSEKVTSALTGVPINTAITSKEISLGYRDYYRTGINMVT  113
AmetDRAFT_1908_1(11-115)       TDQAYGCAQEIRTFKLHLIKRHLVGKIEIGL-PPTFLNHMVNEV  105
AmetDRAFT_1908_2(133-245)      NYQVEGELLLFKKFLRELEALELNQKVLGTL-SALMLDHMAREE  113

consensus/80%                  .tps..tstphhsFKpthhthl.pCcl.u...hPLLsDHs.REA
```

FIGURE 5: BAS0955 is homologous to proteins BT9727_0941 from *Bacillus thuringiensis* serovar konkukian str. 97-27, BCZK0933 from *Bacillus cereus* E33L, and BCE_G9241_1042 from *Bacillus cereus* G9241. BA1021 is homologous to protein GBAA1021 from *Bacillus anthracis* str. "*Ames* Ancestor." BA0807 is homologous to proteins GBAA0807 from *Bacillus anthracis* str. "*Ames* Ancestor" and BAS0770 from *Bacillus anthracis* str. Sterne.

```
Secondary Structure           E E                    HHHHHHHHHH

BT9727_4219_1(67-169)         VDKNGYVVPQTLAIPTPKANE - - - - VIQQTLEYLVKDGPVTNLLPN - GFRAVIPANTSMT - - LDLKKDG
BCZK4235_1(67-169)            VDKNGYVVPQTLAIPTPKANE - - - - VIQQTLEYLVKDGPVTNLLPN - GFRAVIPANTSMT - - LDLKKDG
BA4716_1(67-169)              VDKNGYVVPQTLAIPTPKANE - - - - VIQQTLEYLVKDGPVTNLLPN - GFRAVIPANTSMT - - LDLKKDG
BCE4587_1(67-169)             VDKNGYVVPQTLAIPTPKANE - - - - VIQQTLEYLVKDGPVTNLLPN - GFRAVIPANTSMT - - LNLKKDG
RBTH_05210_1(67-169)          VDKNGYVVPQTLAIPTPKANE - - - - TVKQTLEYLVKDGPVTNLLPN - GFRAVIPANTTMT - - LDLKKDG
ABF83609_1(67-169)            VDKNGYVVPQTLAIPTPKANE - - - - TVKQTLEYLVKDGPVTNLLPN - GFRAVIPANTTMT - - LDLKKDG
BC4495_1(67-169)              VDKNGYVVPQTLAIPTPKANE - - - - TVKQTLEYLVKDGPVTNLLPN - GFRAVIPANTTMT - - LDLKKDG
BcerKBAB4DRAFT_4089_1(67-169) VDKNGYVVPQTLAMPTPKANE - - - - VVQQTLEYLVKDGPVTNLLPN - GFRAVLPANTTMT - - LNLKKGG
Bcer98DRAFT_3179_1(67-169)    VDKNGYVVPQTLPIPKQSE - - - - VVKQTLEYLVKDGPVENILPN - GFRAVLPADTTMT - - VDLKKDG
gerM(84-184)                  IDKNGYVVAQTLPLPKSES - - - - - TAKQALEYLVQGGPVSEILPN - GFRAVLPADTTVN - - VDIKKDG
BSU28380_1(84-184)            IDKNGYVVAQTLPLPKSES - - - - - TAKQALEYLVQGGPVSEILPN - GFRAVLPADTTVN - - VDIKKDG
BL00314_1(87-187)             IDKNGYVTAQTLPLPKQEG - - - - - TAKQALEYLVEGGPVSNILPN - GFRAVLPADTTVN - - VDIKEDG
GK2667_1(76-177)              IDKNGFVVPQTVELPKTQA - - - - - VAKQVLEYLVEDGPVSEMLPN - GFRAVIPAGTTVL - GTKLEKDG
BH3070_1(87-186)              LDENGMVVPQTLPLPKSDG - - - - - VLKQSLEYLVKGGPVTNLLPN - GFQAVLPPDTEMS - - VNL - EDG
OB2107_1(69-172)              LDANGMVASQTLELPVPDTNE - - - - VAAQVLEHLVKGGPVTPLLPN - GFQAVLPEGTEVL - GVNLQEDG
B14911_06091_1(82-181)        VDKNGYVVPQTLTLPKTES - - - - - VATQALEYLMQNGPVTDMLPN - DFRAVLPADTKIS - VN - - VKDK
ABC2653_1(99-200)             IDSNGLVVPQTLTLPKTDS - - - - - VMKQALEYLVEGGPINDILPN - GFRAVLPAGTEVD - IDHLKEEK
SwolDRAFT_2302(77-173)        ADKEELVMERR - EITRTEG - - - - - IARSTLQELLK - GPDN - - - P - - AYRNVFPEGTRLL - DINLKPDG
Moth_0516(72-172)             DSSGNYLVAEKRSIPAVEG - - - - - - IARATIEELIKGPAPDSK - - - - - LLPTIPKGTVLK - DINIRPDG
MothDRAFT_0979(72-172)        DSSGNYLVAEKRSIPAVEG - - - - - - IARATIEELIKGPAPDSK - - - - - LLPTIPKGTVLK - DINIRPDG
CtheDRAFT_0840(63-168)        NEDNSKLKLEIRYIPVSETTKSVNHLAE IIVNELIKGPKVAG - - - - - - LKPTIPEGTKLRSAIKIEGD -
AmetTDRAFT_1640_1(62-164)     RDDKGLLIPVMRRIPWQEG - - - - - IAKAALEQLVDQPVLRDDLATIGLLPVLPPGTEVI - GISINEG -
GAA01614(67-167)              TGSDAYLVREVHQVPFTRE - - - - - VAKAALEELINTAPSTPG - - - - - AVRVLPPATKIR - GISIKDG -
BCE4587_2(220-319)            NNKQQYYVPVTRRVVEGKE - - - - - NDYAAIVDELVKGPIHQS - - - - - - LLNDFNPGVKLI - TNPKLQDG
BA4716_2(220-319)             NNKQQYYVPVTRRVVEGKE - - - - - NDYAAIVDELVKGPIHQS - - - - - - LLNDFNPGVKLI - TNPKLQDG
BT9727_4219_2(220-319)        NNKQQYYVPVTRRVVEGKE - - - - - NDYAAIVDELVKGPIHQS - - - - - - LLNDFNPGVKLI - TNPKLQDG
BCZK4235_2(220-319)           NNKQQYYVPVTRRVVEGKE - - - - - NDYAAIVDELVKGPIHQS - - - - - - LLNDFNPGVKLI - TNPKLQDG
ABF83609_2(20-319)            NNKQQYYVPVTRRVAEGKE - - - - - NDYAAIIDELVKGPIHQS - - - - - - LLNDFNPGVKLI - TNPKLQDG
BC4495_2(220-319)             NNKQQYYVPVTRRVAEGKE - - - - - NDYATIIDELVKGPIHQS - - - - - - LLNDFNPGVKLI - TNPKLQDG
RBTH_05210_2(220-310)         NNKQQYYVPVTRRVAEGKE - - - - - NDYAAIIDELVKGPIHQS - - - - - - LLNDFNPGVKLI - TNPKLQDG
BcerKBAB4DRAFT_4089_2(220-319) NNKQQYYVPVTRRVAEGKE - - - - - NDYSAIVDELVKGPIQGS - - - - - - LLNDFNPGAKLI - TNPKVENG
Bcer98DRAFT_3179_2(219-318)   NNKRQYYVPVTRRVAEGKE - - - - - NEVETIINELVKGPSHYS - - - - - - LLNDFNPGVKLV - SEPKIQDG
BSU28380_2(234-336)           NEDSEYYVPVTKRIDNSEK - - - - - DDITAAINELAKGPSKVSG - - - - - LLTDFSEDVKLV - SKPKIKDG
BL00314_2(237-339)            SDKGTYYVPVTKRTSAKEK - - - - - DQVTAAIKELTEGPDNKSG - - - - - LLSDFQGDVKLE - NKPKIEDG
GK2667_2(227-327)             QGNSTYYVPVTRRVSNKEK - - - - - DDIAAAVNELIQGPEQGSG - - - - - LVGVFQPDAKLV - DAPKYEDG
B14911_06091_2(231-331)       EEGAYYYVPVTKRISAQED - - - - - NQVEAVVKELVKGPSFTSN - - - - - LFTDFMPEVELL - GDPKIENG
BH3070_2(236-335)             SGDQTYYVPVTRRVNVKD - - - - - - NSFATAVEELLNGPMVTSP - - - - - LVTDFRNGVELL - DEPKYENG
ABC2653_2(250-349)            NDEDTYYVPVTKRVENVD - - - - - - NELEAAINELIDGPSLMTN - - - - - LLTEMSGDVELL - NEPKLQNG
OB2107_2(222-322)             QENNRYYVPTQYIETNED - - - - - EAIANIIKELIDGPHQSK - - - - - VVNVFNPEAGLA - SEPTLNNG
AmetDRAFT_1640_2(210-309)     NGEDDFFIPITRGLNVLKA - - - - - - DTKSVLTALVEGAPVGSG - - - - - LHSEIPYGASIN - - DVYVRDG

consensus/80%                 . scptYhVs . Thtlstsct . . . . . . . htthlc . Llcss . hps . . . . . . hhsshssssphh . . . shhp - G
```

```
Secondary Structure           EEEEEE                    HHHHHHHHHH           EEEE

BT9727_4219_1(67-169)         TAVIDFSKEMKNYA - - - - KEEERQIVESIAWTLTQFK - EVKQVQFQ      103
BCZK4235_1(67-169)            TAVIDFSKEMKNYA - - - - KEEERQIVESIAWTLTQFK - EVKQVQFQ      103
BA4716_1(67-169)              TAVIDFSKEMKNYA - - - - KEEERQIVESIAWTLTQFK - EVKQVQFQ      103
BCE4587_1(67-169)             TAVIDFSKEMKNYA - - - - KEEERQIVESIAWTLTQFK - EIKQVQFQ      103
RBTH_05210_1(67-169)          TAVIDFSKEMKNYA - - - - KEEERQIVESIAWTLTQFT - EIKQVQFQ      103
ABF83609_1(67-169)            TAVIDFSKEMKNYA - - - - KEEERQIVESIAWTLTQFT - EIKQVQFQ      103
BC4495_1(67-169)              TAVIDFSKEMKNYA - - - - KEEERQIVESIAWTLTQFT - EIKQVQFQ      103
BcerKBAB4DRAFT_4089_1(67-169) TAVIDFSKEMKNYS - - - - KEEERQIVESVAWTLTQFT - EIKQVQFQ      103
Bcer98DRAFT_3179_1(67-169)    TAVIDFSKEMQNYK - - - - KEEERQIVESVAWTLTQFK - DIKQVKFQ      103
gerM(84-184)                  TAIADFSNEFKNYK - - - - KEDEQKIVQSVTWTLTQFS - SIDKVKLR      101
BSU28380_1(84-184)            TAIADFSNEFKNYK - - - - KEDEQKIVQSVTWTLTQFS - SIDKVKLR      101
BL00314_1(87-187)             TAIADFSNEFKNYK - - - - AEDEQKIVQAITWTLTQFN - SIDKVKLR      101
GK2667_1(76-177)              TLIADFSPEFKNYK - - - - PEDEKRILQSITWTLTQFD - NIKRVKIR      102
BH3070_1(87-186)              VAVVDFSKEFTEYD - - - - GEKEQQIILQSITWTLTQFE - NVEKVKLQ      100
OB2107_1(69-172)              TIIVDLSEEFTQYE - - - - ENQEVQIILESVTHTLTQFE - SVHKVKLR      104
B14911_06091_1(82-181)        VATVDFSKEFGDYQ - - - - AEDEEKILESITWTLTQFD - SIKQVKLR      100
ABC2653_1(99-200)             LAIVNFSSEFNDYN - - - - LADEKQIFEAVTWTLTQFP - DVEEVKVE      102
SwolDRAFT_2302(77-173)        TCILDFSSELRRLEN - - - EVEEKQMLDAVCQTLAQFP - AVKQLVFM      97
Moth_0516(72-172)             LARVDFSKELVANHS - GGSLGESLTVYSIVNTLTQFP - TIKQVQFL      101
MothDRAFT_0979(72-172)        LARVDFSKELVANHS - GGSLGESLTVYSIVNTLTQFP - TIKQVQFL      101
CtheDRAFT_0840(63-168)        VAIVDFTKEFRDNHP - GGKAEERMTIYSVVNSLTELK - EINKVKFL      106
AmetDRAFT_1640_1(62-164)      LSKVDFNEQLLAYQS - - - EIDENAIKSIVYTLTEFD - SIDQVQIM      103
GAA01614(67-167)              LATVDFSRDVLRANT - G - ASGEALGIQSIVNTLTEFP - EVQKVSFL      99
BCE4587_2(220-319)            NLTLNFNENIFINP - DKNMISNYVLKSLVLSLTEKK - GVKSVSIE      100
BA4716_2(220-319)             NLTLNFNENIFINP - DKNMISNYVLKSLVLSLTEKK - GVKSVSIE      100
BT9727_4219_2(220-319)        NLTLNFNENIFINP - DKNMISNYVLKSLVLSLTEKK - GVKSVSIE      100
BCZK4235_2(220-319)           NLTLNFNENIFINP - DKNMISNYVLKSLVLSLTEKK - GVKSVSIE      100
ABF83609_2(20-319)            NLTLNFNENIFVNP - DKNMISNYVLKSLVLSLTEKK - GVKNISIE      100
BC4495_2(220-319)             NLTLNFNENIFVNP - DKNMISNYVLKSLVLSLTEKK - GVKNISIE      100
RBTH_05210_2(220-310)         NLTLNFNENIFVNP - DKNMISNYVLKSLVLSLTEKK - GVKNVSIE      100
BcerKBAB4DRAFT_4089_2(220-319) NITLNFNENIFVNP - DKNMISNYVLKSLVLSLTEKQ - GVKNVSIE      100
Bcer98DRAFT_3179_2(219-318)   KVTLNFNENIYANK - DKNMISNYVLQSLVLSLTEKQ - GVKNVSVE      100
BSU28380_2(234-336)           RVTLDFNQSIFGSADEKTKMISSEVLNSIVLTLTEQP - DVKSVSVK      103
BL00314_2(237-339)            HVTLDFNEAIYGSADGQKKVISDVVLNSIVLTLTELP - DVKSVSVT      103
GK2667_2(227-327)             KVTLNFNEGIYGSN - - KKNVISDVVLNSLVLSLTEQK - GVESVAIT      101
B14911_06091_2(231-331)       LATLDFNESVYGSF - - EEKIISQHLLNSLVLSLTEQK - GIESVAVT      101
BH3070_2(236-335)             VVTLNFNEALLSQM - - QATAVSDEIINMLALTLTEQD - GVEKVAIQ      100
ABC2653_2(250-349)            EVVLDFNEAIQSAN - EGSAIPTSVLESLALTLTEQG - GIEKVSIQ      100
OB2107_2(222-322)             ILEVVFNKEILLADS - EQGIIADEVMETMVRTLTEQP - NIDAVDVK      101
AmetDRAFT_1640_2(210-309)     IAYIDFTEEIRNVP - VNEKHQQSLVYELGLTLREVEPSIHQVRIL      100

consensus/80%                 . hhlsFscphhs . . . . . . . p . . pp . llpSlshoLTpht . tlcpVphp
```
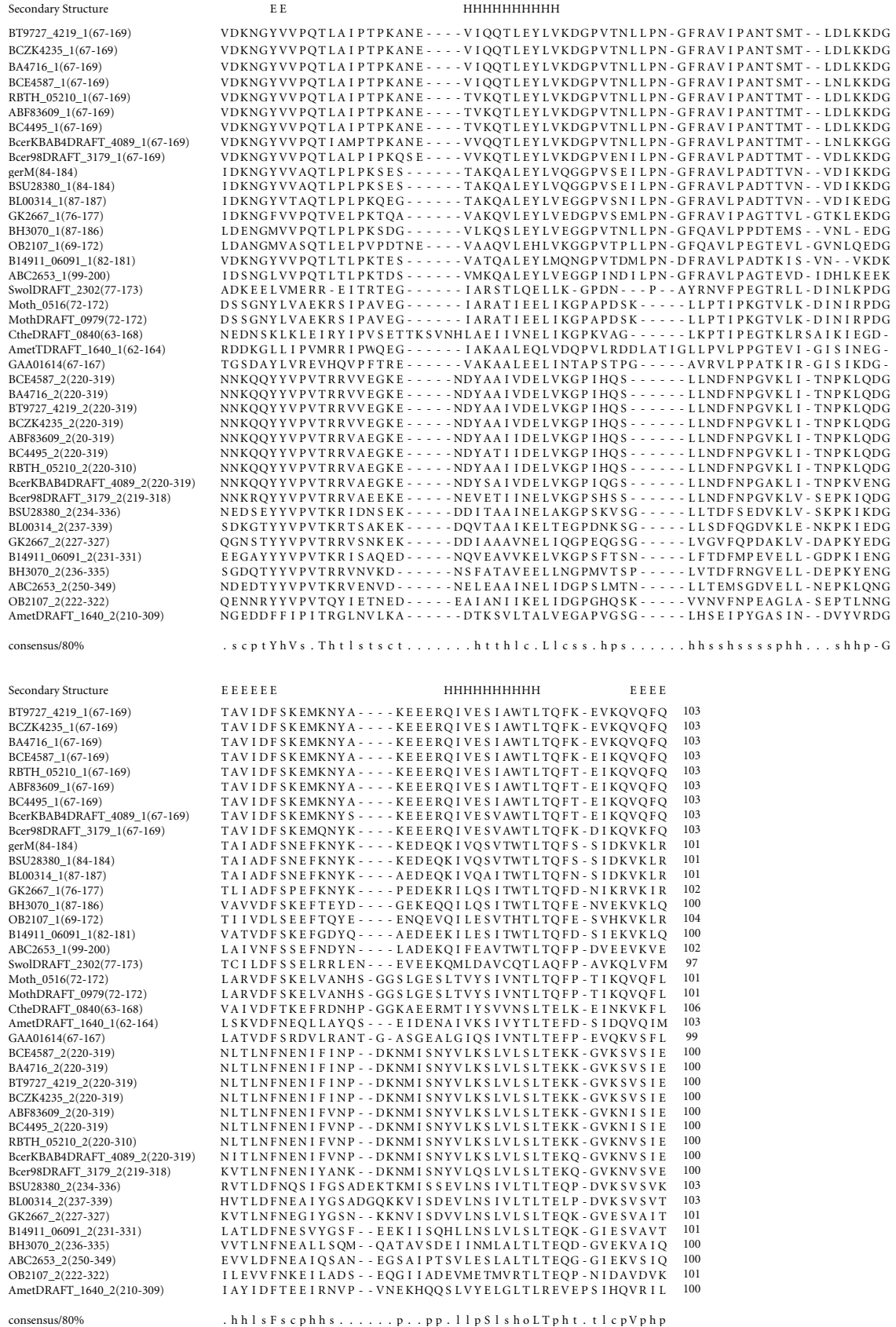
FIGURE 6: BA4716 is homologous to proteins GBAA4716 from *Bacillus anthracis* str. "*Ames* Ancestor," BAS4378 from *Bacillus anthracis* str. Sterne, and Bant_01005366 from *Bacillus anthracis* str. A2012. BT9727_4219 is homologous to protein BCZK4235 from *Bacillus cereus* E33L. BA4716 is homologous to protein BL02986 from *Bacillus licheniformis* ATCC 14580.

```
Secondary structure                          E E E E E                       E E E E E                      E E E E E
BC4088_1(47-130)                  I K P G E K T E V Q A L V T Q G K E K V T D A D D V K F E I W K D G D - - E K H E M L D G K H K G K G V Y A V E K T F E T D G
RBTH_02670_1(47-130)              I K P G E K T E V Q A L V T Q G K E K X T D A D D V K F E I W K D G D - - E K H E M L D G K H K G K G V Y A V E K T F E T D G
BCE_G9241_4093_1(45-128)          I K P G E K T E V Q A L V T Q G K E R V T D A D D V K F E I W K D G D - - E K H E M L D G K H K G K G V Y A V E K T F E T D G
BA4310_1(45-128)                  I K P G E K T E V Q A L V T Q G K E K V T D A D D V K F E V W K A G D - - E K H E M L E G K H K G K G V Y A V E K T F E T D G
BT9727_3829_1(45-128)             I K P G E K T E V Q A L V T Q G K E K V T D A D D V K F E V W K A G D - - E K H E M L E G K H K G K G V Y A V E K T F E T D G
Bant_01004966_1(51-134)           I K P G E K T E V Q A L V T Q G K E K V T D A D D V K F E V W K A G D - - E K H E M L E G K H K G K G V Y A V E K T F E T D G
BCE4157_1(45-128)                 I K P G E K T E V Q A L V T Q G K E K V T D A D D V K F E I W K A G D - - E K H E M L E G K H K G K G V Y A V E K T F E T D G
BCZK3845_1(45-128)                I K P G E K T E V Q A L V T Q G K E K V T D A D D V K F E I W K A G D - - E K H E M L E G K H K G K G V Y A V E K T F E T D G
BcerKBAB4DRAFT_2040_1(46-128)     I K P G E K T E V Q A L V T Q G K E K V T D A D D V K F E I W K A G D - - E K H E M L N A K H K G K G V Y A V E K T F E T D G
GK0969(45-128)                    I D L N K P T K L A C V V T Y G G E K D E A N E V K F E V W K H G S - - D E R E M L E A K H D G D G R Y S V E K T F T E A G
BL05305(45-129)                   A A K N E K A V I K A T V L Y G E E P V A D A D E V E F E C W K A G S K - E D S E L I K A K N E G K G V Y S M E K A F P E D G
BSU30660(44-127)                  V N P G E S A A Y E A A V S Y G D E A V T D A D E V E F E V W K E G E K - D A S Q M F K V K Q E - K G V Y R L E T T F K E D G
OB2488(50-134)                    V E T G E T I D L T A H V T Y G D A P V E D A D E V I F E V W T Q G N S - D Q S V E L E G K H Q E N G T Y T A S Y T F E E E K
B14911_05359_1(53-137)            V E L N E E I T L S V E V V Q G E E A V E D A D E V K F E I W Q E G N Q - E E S E M L P A E H T G K G I Y Q A A K T F G K D G
BH0678_1(45-129)                  L A S G E N M T F D V L V T Q N E A P V E D A R E V I V E F W Q E G A K - E E S D M I E S T N E G G G V Y R V T Y E F P E D G
ABC0230(45-129)                   I E I G E E I L L S V Q L A Q G E V Q V E D A D E V V F E V W K D Q E R - D N G T L Q E A T H Q E N G V Y E I T H T F D E D G
ABC4088(44-127)                   L E L - E N I V L E A K V M Q G D E P V D D A E E V V F E V W P Y D D R - E E S E F H E A S Y A E S G L Y Q A P L A L E E A G
BH0983(47-131)                    L I P N T P H E L A I H V T Q G D E N V T D A T D I Q F E I W Q G H D R - E Q G E L I E A S H V E D G I Y L V E Y E F P E D G
B14911_09907(34-118)              F A A G E D V P I R A V L T Q N G E K V A G A D Y V H F E I W K R D G S - V H Y P M E E A A D E G E G V Y Q L T K K F E Q D G
ExigDRAFT_1796(51-135)            A D Q E K Q Y R F G A T L W Q D Q K A V K E A E Y V H F E I W K A D G T - L R Y S M E P A D E T K P G V Y S I E K K L P K E G
BAA83944_1(46-130)                L V T D Q E E S L T V S L S H N G E I L S K V D S L H V H I W K H D H T - V A Y H F E Q L E T D Q D G A F N L P L T F E S D G
BH1853(46-130)                    L V T D Q E E S L T V S L S H N G E I L S K V D S L H V H I W K H D H T - V A Y H F E Q L E T D Q D G A F N L P L T F E S D G
OB3282(48-131)                    I E A K E N T E V T F E L S Q N G E S V S T L D D L S V T T W M V D S E - T T K Q L V A E N V G - N G E Y S V E T S F D Q D G
BCE_G9241_4093_2(163-245)         I K A N A E S T M K V H L K Q K E - E A L T G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E T K H T F K E N G
BC4088_2(165-247)                 I K A N A E S T M K V H L K Q K E - E A L T G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E T K H T F K E N G
RBTH_02670_2(165-247)             I K A N A E S T M K V H L K Q K E - E A L A G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E T K H T F K E N G
BcerKBAB4DRAFT_2040_2(158-240)    I K A N A E S T M K V H L K Q K E - E A L S G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E S K H T F K E N G
BT9727_3829_2(163-245)            I K A N A E S T M K V H L K Q K E - E A L T G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E T K H T F K E N G
Bant_01004966_2(169-251)          I K A N A E S T M K V H L K Q K E - E A L T G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E T K H T F K E N G
BA4310_2(163-245)                 I K A N A E S T M K V H L K Q K E - E A L T G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E T K H T F K E N G
BCE4157_2(163-245)                I K A N A E S T M K V H L K Q K E - E A L T G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E T K H T F K E N G
BCZK3845_2(163-245)               I K A N A E S T M K V H L K Q K E - E A L T G A E V Q L E I W K D G V - - E K H E F I P A K E G N K G E Y E T K Y T F K E K G
Bcer98DRAFT_3614(94-176)          V K A N A E S T L K A H V K Q K E - E A L T K A E V Q F E I W K D G V - - E K H T F I T A K E D N K G E Y V G K Y T F K E S G
B14911_05359_2(187-271)           I H M K Q A A G L D V Q V D K K D G A P L E K A L V K L E I M K E G K - - D T P E W V N L K E S G E G K Y S A E H S F A E A G
BH0678_2(159-242)                 I Q A G E E T T L L I V V E H K D - K P F T G G V T L E V W Q H E D - - E A H T W L D T E E T D V G Q Y E V S H T F A D A G
ExigDRAFT_0574(52-137)            K T M E N Q K V V F Q A T A L E N K K A V N L E N V A F E V W K A D E K E A V H Q K F K A A L K K T G T Y Q A E A K L A - E G

consensus/80%                     l p . s t p t p h p s h l p p t c . t s . s u s - V p h E l W K t s s . . - p p p h h . u c p t t p G . Y t s p h o F t p s G
```

```
Secondary structure                          E E E E E E E                 E E
BC4088_1(47-130)                  V Y H I I A H T N A R E - M H V M P E V K V A V   84
RBTH_02670_1(47-130)              V Y H I I A H T N A R E - M H V M P E V K V A V   84
BCE_G9241_4093_1(45-128)          V Y H I I P H T N A R D - M H V M P E H K V A V   84
BA4310_1(45-128)                  V Y H I I A H T N A R E - M H V M P E V K V A V   84
BT9727_3829_1(45-128)             V Y H I I A H T N A R E - M H V M P E V K V A V   84
Bant_01004966_1(51-134)           V Y H I I A H T N A R E - M H V M P E V K V A V   84
BCE4157_1(45-128)                 V Y H I I A H T N A R E - M H V M P E V K V A V   84
BCZK3845_1(45-128)                V Y H I I A H T N A R E - M H V M P E V K V A V   84
BcerKBAB4DRAFT_2040_1(46-128)     V Y H V I A H T N A R E - M H V M P E V K V A V   84
GK0969(45-128)                    T Y S V V A H V T A R D - M H N M P K K D I V A   84
BL05305(45-129)                   H Y K V Q V H V T A K K - Q H T M P V A D I K V   85
BSU30660(44-127)                  V Y T V Q S H V T A K K - Q H S M P T L K V Q V   84
OB2488(50-134)                    V Y E M Y A H T T A E A - I H S M P F K T V I V   85
B14911_05359_1(53-137)            D Y I V Q V H V T A R D - M H T M P K A E V Q A   85
BH0678_1(45-129)                  L Y F V Q P H V T A R D - M H R M P L Y E L T I   85
ABC0230(45-129)                   I Y I V Q T H V T A R D - M H V M P K Q M I V A   85
ABC4088(44-127)                   I Y M V Q V H V T A R G - M H V M P T Q P L F A   84
BH0983(47-131)                    I Y F V Q A H V T A R G - L H V M P T E R L I V   85
B14911_09907(34-118)              V Y I I K V H A S S G G - S L I M P Q V K A V   85
ExigDRAFT_1796(51-135)            L Y Y I K V H A S S N G - A M I M P T R Q F I V   85
BAA83944_1(46-130)                L Y Y M K V D V T H N G - D T I M P T A Q L I V   85
BH1853(46-130)                    L Y Y M K V D V T H N G - D T I M P T A Q L I V   85
OB3282(48-131)                    I Y H M K V T A S K N N - A T I M P T K Q F I V   84
BCE_G9241_4093_2(163-245)         A Y K V K V H V R K G E - L H E H K E E T I E V   83
BC4088_2(165-247)                 A Y K V K V H V R K G E - L H E H K E E T I E V   83
RBTH_02670_2(165-247)             A Y K V K V H V R K G E - L H E H K E E T I E V   83
BcerKBAB4DRAFT_2040_2(158-240)    A Y K V K V H V R K G E - L H E H K E E T V E V   83
BT9727_3829_2(163-245)            S Y K V K V H V K K G E - L H E H K E E T V E V   83
Bant_01004966_2(169-251)          S Y K V K V H V K K G E - L H E H K E E T V E V   83
BA4310_2(163-245)                 S Y K V K V H V K K G E - L H E H K E E T V E V   83
BCE4157_2(163-245)                S Y K V K V H V K K G E - L H E H K E E T V E V   83
BCZK3845_2(163-245)               S Y K V K V H V K K G E - L H E H K E E T V E V   83
Bcer98DRAFT_3614(94-176)          K Y K V K V H V R K G D - L H E H K E E T V E V   83
B14911_05359_2(187-271)           S Y T V T V H V E N S E G L H E H S D F P L T V   85
BH0678_2(159-242)                 E Y H V V F H I E D D T G L H E H I H E A L I V   84
ExigDRAFT_0574(52-137)            E Y E G L Y H I N D K N G L H H M D K I S F V V   86

consensus/80%                     . Y h l h s H s p t t p . h H . h . p . p l . V
```

Figure 7: BA4310 is homologous to proteins GBAA4310 from *Bacillus anthracis* str. "*Ames* Ancestor," BAS3998 from *Bacillus anthracis* str. Sterne, and BT9727_3829 from *Bacillus thuringiensis* serovar konkukian str. 97-27.

Secondary structure                        E E E E               H H        H                        H H H H H

BCZK2413_2(120-222)          V Y N T G F I G V V F A D L C S I D R F N F E F - - - E M G M L T K L M K D M I I P V K E L F L R H N V P A Y I S T S H L E E Q N K
BT9727_2444_2(120-222)       V Y N T G F I G V V F A D L C S I D R F N F E F - - - E M G M L T K L M K D M I I P V K E L F L R H N V P A Y I S T S H L E E Q N K
BA2665_2(120-222)            V Y N T G F I G V V F A D L C S I D R F N F E F - - - E M G M L T K L M K D M I I P V K E L F L R H N V P A Y I S T S H L E E Q N K
Bant_01003317_2(124-226)     V Y N T G F I G V V F A D L C S I D R F N F E F - - - E M G M L T K L M K D M I I P V K E L F L R H N V P A Y I S T S H L E E X N K
BCE2700_2(122-224)           V F N T G F I G V V F A D L C S I D R F N F E F - - - E M G M L T K L M K D M I I P V K E L F L R H N V P A Y I S T S H L E E Q N K
BCE_G9241_CNI_0263_2(122-224) V F N T G F I G V V F A D L C S I D R F N F E F - - - E M G M L T K L M K D M I I P V K E L F L R H N V P A Y I S T S H L E E Q N K
BcerKBAB4DRAFT_0535_2(120-222) V F N T G F I G V V F A D L S S I D R F N F E F - - - E M G M L T K L M K D M I I P V K E L F L R H N V P A Y I S T S H L E E Q N K
BC2674_2(122-224)            V F N T G F I G V V F A D L C S I D R F N F E F - - - E M G M L T K L M K D M I I P V K E L F L R H N V P A Y I S T S H L E E Q N K
Bcer98DRAFT_0128_2(122-224)  V F N T G F I G V V F A D L S S I D R F N F E F - - - E M N M L F K L M K D M I I P V K E L F L R H N I P A Y I S T S H L E T Q N K
BA2665_1(16-119)             I S N T G F I G S V F I D T L E L Q K K S Y Y F A R K K L Q I V H H V L D G L S G A T S S L F K E H N I S A Y M S C V Y L H K Q K K
Bant_01003317_1(20-123)      I S N T G F I G S V F I D T L E L Q K K S Y Y F A R K K L Q I V H H V L D G L S G A T S S L F K E H N I S A Y M S C V Y L H K Q K K
BCZK2413_1(16-119)           I S N T G F I G S V F I D T L E L Q K K S Y Y F A R K K L Q I V H H V L D G L S G A T S S L F K E H N I S A Y M S C V Y L H K Q K K
BT9727_2444_1(16-119)        I S N T G F I G S V F I D T L E L Q K K S Y Y F A R K K L Q I V H H V L D G L S G A T S S L F K E H N I S A Y M S C V Y L H K Q K K
BcerKBAB4DRAFT_0535_1(16-119) I S N T G F I G S V F I D T L E L Q K K S Y Y F A R K K L Q I V H H V L D G L S G A T S A L F K E H N T A A Y M S C V Y L H K Q K K
BCE2700_1(16-121)            I S N T G F I G S V F I D T L E L Q K K S Y Y F A R K K L Q I V H H V L D G L S G A T S S L F K E H N I S A Y M S C V Y L H K Q K K
BCE_G9241_CNI_0263_1(16-121) I S N T G F I G S V F I D T L E L Q K K S Y Y F A R K K L Q I V H H V L D G L S G A T S S L F K E H N I S A Y M S C V Y L H K Q K K
BC2674_1(16-121)             I S N T G F I G S V F I D T L E L Q K K S Y Y F A R K K L Q I V H H V L D G L S G A T S A L F K E H N I S A Y M S C V Y L H K Q K K
Bcer98DRAFT_0128_1(16-121)   I S N T G F I G S V F I D T L E L Q K K S Y Y F S R K K L Q I V H H V L D G L A E A T S S L F H E H E V A A Y I S C V Y L H K Q K K

consensus/80%                l . N T G F I G s V F h D h h p l p + h s a . F . . . c h t h l p + l h c s h . h s s p p L F h c H N l s A Y h S s s a L c c Q p K


Secondary structure                        E E E E               H H H H H H H H H H H

BCZK2413_2(120-222)          L G F V L S I K P Y D E R A E A D L Y F E A Y L K E R G L F I G - D E E D D I D K    103
BT9727_2444_2(120-222)       L G F V L S I K P Y D E R A E A D L Y F E A Y L K E R G L F I G - D E E D D I D K    103
BA2665_2(120-222)            L G F V L S I K P Y D E R A E A D L Y F E A Y L K E R G L F I G - D E E D D I D K    103
Bant_01003317_2(124-226)     L G F V L S I K X Y D E R A E A D L Y F E A Y L K E R G L F I G - D E E D D I D K    103
BCE2700_2(122-224)           L G F V L S V K P Y D E R A E A D L Y F E A Y L K E R G L F I G - D E E D D I D K    103
BCE_G9241_CNI_0263_2(122-224) L G F V L S V K P Y D E R A E A D L Y F E A Y L K E R G L F I G - D E E D D I D K    103
BcerKBAB4DRAFT_0535_2(120-222) L G F V L S V K P Y D E R A E A D L Y F E T Y L K E R G L F I G - D E E D D I D K    103
BC2674_2(122-224)            L G F V L S V K P Y D E R A E A D L Y F E A Y L K E R G L F I G - D E E D D I D K    103
Bcer98DRAFT_0128_2(122-224)  V G F V L S I K P Y D E R A E A D L Y F E T Y L K E R G L F I G - D E E D E M D K    103
BA2665_1(16-119)             I G F V L S T K P F E Q - S D G V A Y F I N Y L I E K N F Y G - - N E E V E Y Q E    104
Bant_01003317_1(20-123)      I G F V L S T K P F E Q - S D G V A Y F I N Y L I E K N F Y G - - N E E V E Y Q E    104
BCZK2413_1(16-119)           I G F V L S T K P F E Q - S D G V A Y F I N Y L I E K N F Y G - - N E E V E Y Q E    104
BT9727_2444_1(16-119)        I G F V L S T K P F E Q - S D G V A Y F I N Y L I E K N F Y G - - N E E V E Y Q E    104
BcerKBAB4DRAFT_0535_1(16-119) I G F V L S T K P F E Q - S D G V S Y F I N Y L I E K N F Y G - - N E E V E Y Q E    104
BCE2700_1(16-121)            I G F V L S T K P F E Q - S D G V A Y F V N Y L I E K N F Y G N H D E D V E Y Q E    106
BCE_G9241_CNI_0263_1(16-121) I G F V L S T K P F E Q - S D G V A Y F V N Y L I E K N F Y G G H D E D V E Y Q E    106
BC2674_1(16-121)             I G F V L S T K P F E Q - S D G V A Y F V N Y L I E K N F Y G G H D E D V E Y Q E    106
Bcer98DRAFT_0128_1(16-121)   I G F V L S T K L F E Q - T D G I A Y F K N Y L I E K N F Y G K T D Q E V E Y Q E    106

consensus/80%                l G F V L S h K P a - p . u - u s h Y F . s Y L h E + s h a h . . s E E s - h p c

FIGURE 8: BA2665 is homologous to proteins GBAA2665 from *Bacillus anthracis* str. "*Ames* Ancestor," BAS2482 from *Bacillus anthracis* str. Sterne. BT9727_2444 is homologous to protein BCZK2413 from *Bacillus cereus* E33L.

varied between 11–98%. The secondary structure is predicted to comprise two $\alpha$-helices and three $\beta$-strands as shown in Figure 6. The representative domain architecture corresponding to proteins comprising this domain is shown in Figure 19.

### 3.6. 84-amino-acid-residue ExW domain

The 246-amino-acid-residue protein corresponding to the GENE_ID BA4310 and described as hypothetical protein comprises an 84-amino-acid-residue region as two copies. Further BLAST searches using sequence corresponding to the domain (45–128) as a query identified 25 proteins (Table 1(f)) that are described as either conserved or hypothetical proteins. This domain region occurs as two copies in proteins of *B. anthracis* str. *Ames, B. cereus, B. halodurans* (GENE_ID BH0678), *B. thuringiensis*, *B. thuringiensis serovar israelensis, Geobacillus kaustophilu*s, *Bacillus weihenstephanensis*, and *Exiguobacterium sibiricum* genomes and as single copy in proteins of *B. clausii, B. halodurans* (GENE_ID BH0983), *B. licheniformis, B. subtilis, Exiguobacterium sp.*, and *Oceanobacillus ihenyensis* genomes. The length of proteins varied between 142 to 273-amino-acid residues. The

multiple sequence alignment corresponding to this domain identified ExW sequence motif. The pairwise sequence identities corresponding to the ExW domain varied between 14–98%. The secondary structure of this domain is predicted to comprise five $\beta$-strands and the conserved sequence motif is associated with one of the $\beta$-strands as shown in Figure 7. The representative domain architecture corresponding to proteins comprising this domain is shown in Figure 20.

### 3.7. 104-amino-acid-residue NTGFIG domain

The 232-amino-acid-residue protein corresponding to the GENE_ID BA2665 and described as a hypothetical protein comprises a 104-amino-acid-residue region as two copies in tandem. Further BLAST searches using sequence corresponding to the region (16–119) as query identified 9 hypothetical proteins comprising this domain from organisms such as *B. anthracis, B. thuringiensis, Bacillus weihenstephanensis*, and *B. cereus*. The protein corresponding to the GENE_ID BCZK2413 of *B. cereus* is described as group-specific protein. The list of 9 proteins comprising this domain is shown in Table 1(g). The length of proteins varied between 232 to 236-amino-acid residues. This domain

| Secondary structure | HHHHHHHH | |
|---|---|---|
| BT9727_3378_2(139-176) | YKTTI SAQF EYNR FTRDF F EDPNNKGKS KADA I AAWNE | 38 |
| BcerKBAB4DRAFT_0944_2(139-176) | YKTTI SAQF EYNR FTRDF F EDPNNKGKS KADA I AAWNE | 38 |
| BA3686_2(139-176) | YKTTI SAQF EYNR FTRDF F EDPNNKGKS KADA I AAWNE | 38 |
| BCZK3328_2(139-176) | YKTTI SAQF EYNR FTRDF F EDPNNKGKS KADA I AAWNE | 38 |
| BCE_G9241_3579_2(139-176) | YKTTI SSQF EYNR FTRDF F EDPNNKGKS KADA I AAWNE | 38 |
| BCE3645_2(139-176) | YKTTI SPQF EYNR FTRDF F EDPNNKGKT KADV I AAWNE | 38 |
| RBTH_03615_2(139-176) | YKTTI GAQF EYNR FTRDF F EDPNNKGKS KADA I AAWNE | 38 |
| BC3626_2(139-176) | YKTTI GTQF EYNR FTRDF F EDPNNKGKA KADA I AAWNE | 38 |
| B14911_25780_2(138-175) | YKSE I GRQF EYNQF I RDYYADQKNQGKS RAEA I AAWML | 38 |
| RBTH_03615_1(94-129) | FKEK I GTNF RFTVALQKFFK - - ENVGKTYEDAVA FWHE | 36 |
| BC3626_1(94-129) | FKEK I GTNF RFTVALQKFFK - - ENVGKTYEDA I AFWHE | 36 |
| BT9727_3378_1(94-129) | FKEK I GANF RFTVALQKFFK - - EN I GKTYEDAVA FWHE | 36 |
| BCZK3328_1(94-129) | FKEK I GANF RFTVALQKFFK - - EN I GKTYEDAVA FWHE | 36 |
| BA3686_1(94-129) | FKEK I GANF RFTVALQKFFK - - EN I GKTYEDAVA FWHE | 36 |
| BCE_G9241_3579_1(94-129) | FKEK I GANF RFTVALQKFFK - - ENVGKTYEDA I TFWYE | 36 |
| BCE3645_1(94-129) | FKEK I GANF RFTVALQKFFK - - ENVGKTYEDA I TFWYE | 36 |
| BcerKBAB4DRAFT_0944_1(94-129) | FKEK I GANF RFTVALQKFFK - - ENVGKTYEDA I TFWYE | 36 |
| B14911_25780_1(93-128) | FKSV I GSHFHFSTY IQDYFK - - HNPGKTYNDAVS AWHE | 36 |
| consensus/80% | aKpp I uspFcashhhpcFFc. . pNhGKohtDA l uhW. E | |

FIGURE 9: BT9727_3378 is homologous to protein BCZK3328 from *Bacillus cereus* E33L. BA3686 is homologous to proteins GBAA3686 from *Bacillus anthracis* str. "*Ames* Ancestor," BAS3417 from *Bacillus anthracis* str. Sterne, and Bant_01004341 from *Bacillus anthracis* str. A2012.

| Secondary structure | E E E          E E E E          E E E |
|---|---|
| BAS1577_2(128-220) | AKKYDTQV S LAPAVKN I V I LNND - DADD I VRVTGLE SGDVVKVYGEATGG - EV I EKATVQG |
| BA1701_2(126-218) | AKKYDTQV S LAPAVKN I V I LNND - DADD I VRVTGLE SGDVVKVYGEATGG - EV I EKATVQG |
| RBTH_03882_10(898-990) | AVKYE SQVTAE PVGGN I VVLNND - GAAD I VRVTGLTAGDKV SVYNE ETVQ - EA I GTATVAE |
| BAS1577_1(33-127) | AAEVA I VKTKAVTVDA I TVANNEKEAEDT I KVTGLVTGD I VKVYDAAS KGKELGTTK - VAE |
| BA1701_1(31-125) | AAEVA I VKTKAVTVDA I TVANNEKEAEDT I KVTGLVTGD I VKVYDAAS KGKELGTTK - VAE |
| RBTH_03882_2(610-705) | VKYEAE PTTVAPAVEK I TV SNNKVEAEDT I TV SELKKGD I VRVYEAS KGGEA I VTS EAVAE |
| RBTH_03882_3(802-897) | VKYEAE PTTVAPAVEK I TV SNNKVEAEDT I TV SELKKGD I VRVYEAS KGGEA I VTS EAVAE |
| RBTH_03882_1(418-513) | VKYEAE PTTVAPAVEK I TV SNNKVEAEDT I TV SELKKGD I VRVYEAS KGGEA I VTS EAVAE |
| RBTH_03882_8(226-320) | VKYEAE PTTVAPAVEK I TV SNNKVGNADA I TV SKLKKGD I VRVYEAS KGGAA I AAS EAVAE |
| RBTH_03882_6(33-128) | AAEVT SAKTAAL SVEKAN I I NNKKGETDT I TV SELKKGD I VRVYEAS KGGEA I ATS EAVAE |
| RBTH_03882_4(321-416) | AVKYE SQVTVAPAVDTVKVANNKAGDADT I TV SGVAEGDLVRVYDASTEG - KELGNATVAK |
| RBTH_03882_5(706-800) | AVKYE SQVTVAPAVDTVKVANNKAGDADT I TV SEVTEGDVVKVYDASTEG - KELGNATVAK |
| RBTH_03882_7(514-608) | AVKYE SQVTVAPAVDTVKVANNKAGDADT I TV SGVAEGDLVRVYDASTEG - KELGNATVAK |
| RBTH_03882_9(129-224) | AMKYE SEVTVAPAVDTVKVANNKAGDADT I TV SELAPGD I VK I YDASTGGNLKATSAAVAE |
| DSY3134_1(51-142) | VP FSEPLKTTTP - - SA I EVRNY I EG I RDRVTVS SLEEGD I VK I YPSEESN - TPSGTEAVKA |
| DSY3134_2(150-240) | P I PWL I YGHTGNWGEDVKLP RTPFDQS K - AS YPAY - P I DANG I SDDNPLG I I YNQH I I I KG |
| consensus/80% | sh. . t . . hThAssVcp l p l . NNc . t st Dh l pVot l t pGD l V+VYp supt G. t . hssts Vt t |

| Secondary structure | E E E          E E E E E E | |
|---|---|---|
| BAS1577_2(128-220) | NKTAVNVK I PQLG I EAG - KVYVTVTKPNKDE SKRV | 93 |
| BA1701_2(126-218) | NKTAVNVK I PQLG I EAG - KVYVTVTKPNKDE SKRV | 93 |
| RBTH_03882_10(898-990) | NKTAVNVV I PQLGEVAG - K I YVSVTKVNKDE SKRV | 93 |
| BAS1577_1(33-127) | NATDAT I TGKDLLAVAGGTVYVSVQS KDQLE SPRT | 95 |
| BA1701_1(31-125) | NATDAT I TGKDLLAVAGGTVYVSVQS KDQLE SPRT | 95 |
| RBTH_03882_2(610-705) | GKTEAT I LGKDLLKVTGGTVYVSVQS ENELE SART | 96 |
| RBTH_03882_3(802-897) | GKTEAT I LGKDLLKVTGGTVYVSVQS ENELE SART | 96 |
| RBTH_03882_1(418-513) | GKTEAT I LGKDLLKVTGGTVYVSVQS ENELE SART | 96 |
| RBTH_03882_8(226-320) | GKTEAT I LGKDLLKVTGGTVYVSVQS ENELE SART | 96 |
| RBTH_03882_6(33-128) | GKVEVT I TKKDLLKATGGTVYVSVQS E SELE STRT | 96 |
| RBTH_03882_4(321-416) | DAKEAT I TGKDLLV STGGTVYVTVTKPNKDE SKRV | 95 |
| RBTH_03882_5(706-800) | DAKEAT I TGKDLLV STGGTVYVTVTKPNKDE SKRV | 95 |
| RBTH_03882_7(514-608) | EATEVK I EKTDLLV STGGTVYVTVTKPNKDE SKRV | 95 |
| RBTH_03882_9(129-224) | GKKEAT I TGKDLLV STGGTVYVTVTKPNKDE SKRV | 96 |
| DSY3134_1(51-142) | GQTSVT I E I DQLSEVYG - E I YVTVTR SGYEE SDRV | 92 |
| DSY3134_2(150-240) | NGSRVTFYG - - - YAQNAYKDF I LLPSE SVAKKT I E | 91 |
| consensus/80% | st sp ss l hh . pLhhssG. pVYVoVpp . sp . E St Rs | |

FIGURE 10: BA1701 is homologous to proteins GBAA1701 from *Bacillus anthracis* str. "*Ames* Ancestor," and Bant_01002313 from *Bacillus anthracis* str. A2012.

occurs twice in every protein of the bacillus species as shown in Table 1(g). We refer to this as the NTGFIG domain based on the conserved sequence motif that is present at the N-terminal part. The pairwise sequence identities between sequences corresponding to this domain varied between 31–99%. The secondary structure corresponding to this domain is predicted to comprise three $\alpha$-helices and two $\beta$-strands as shown in Figure 8. The representative domain architecture corresponding to proteins comprising this domain is shown in Figure 21.

### 3.8. 36-amino-acid-residue NxGK repeat

The 193-amino-acid-residue protein corresponding to GENE_ID BA3686 and described as hypothetical cytosolic protein comprises a 36-amino-acid-residue region as two

| Secondary structure | HHHHHHHHHHHHHHHHH | HH | HHHHHH |
|---|---|---|---|
| RBTH_06405_4(259-331) | R E K A L D A L Q W T I E E K E K L T D N Q L L Q Q Y T M Q W L K N H R L W T P V V R Y W N G S P Y A M I N D L Y P N K Y |
| pBMB165_3(175-247) | R E K A L E A L Q W T I E E K E K L I D N Q L L Q Q Y T M K W L K R H R L W T P V V R Y W N G S P Y A M I N D L Y P N K Y |
| pE33L466_0092_4(259-328) | K R K A L E A L R W T I E E K E K L D E K Q L L K V F N Q K W L I K Q K L W T P L K R Y W K G S P Y E M L I A L Y P N R F |
| RBTH_06405_3(109-183) | K E K A L Q L L K W L I E E E E K L P P Q K L L Q I Y G Q K W L I E H R L S A P L R V I W N G S P Y A M I N D L Y P N R F |
| pBMB165_2(25-99) | K E K A L Q L L K W I I E E E E K V S P Q K L L Q I Y G Q K W L N E R R L S A P L R V I W D G S P Y A M I N D L Y P N R F |
| BA3147_2(109-183) | K E K A L E A L K W T V E E K E K L S K V E L L K F Y S K K W L E K N K L S A P L V M Y W N G S P Y A M I N S L Y P N K F |
| Bant_01003795_1(25-99) | K E K A L E A L K W T V E E K E K L S K V E L L K F Y S K K W L E K N K L S A P L V M Y W N G S P Y A M I N S L Y P N K F |
| BAS2924_2(116-190) | K E K A L E A L K W T V E E K E K L S K V E L L K F Y S K K W L E K N K L S A P L V M Y W N G S P Y A M I N S L Y P N K F |
| BAS2924_3(191-265) | K E K A L E A L K W T V E E K E K L S K V E L L K F Y S K K W L E K N K L S A P L V M Y W N G S P Y A M I N S L Y P N K F |
| pE33L466_0092_2(109-183) | K E K A L T I L K W I I E E K E G L S Q E K L L E L Y G K K W L E K N K L G A P L A M Y W N S S P Y A M I N D L Y P R R F |
| RBTH_06405_2(184-258) | K D K T L Q A L K W T I E K K E K L N V D Q L K N I Y D N K W L V Q S G L S G A C Q L Y W N D S P Y A M I N D L Y P G Q F |
| pBMB165_1(100-174) | K E K A L Q A L K W T I E E K E K L N P D Q L K N I Y E N K W L T Q L G L R G A C Q L Y W N D S P Y A M I N D L Y P N Q F |
| BAS2924_4(266-340) | K E K A L V A L R W T I E E K E K L T S F Q L L Q V Y S V K W L T I H N L I S P C Q I F W N N S P Y S M I N E L Y P G Q N |
| Bant_01003795_2(100-174) | K E K A L V A L R W T I E E K E K L T S F Q L L Q V Y S V K W L T I H N L I S P C Q I F W N N S P Y S M I N E L Y P G Q N |
| BA3147_3(184-258) | K E K A L V A L R W T I E E K E K L T S F Q L L Q V Y S V K W L T I H N L I S P C Q I F W N N S P Y S M I N E L Y P G Q N |
| pE33L466_0092_3(184-258) | K E K A L E A L K W T I E E K E G L T P K Q L L D V Y N I K W L Q T H R L A S A C Q I I W G N S P F R M I N D L Y I D R F |
| BA3147_1(34-108) | R E L S K R V T K Y L I E T I L K W N E E D I K Q K W N T P L I I K Y R L L G A L K H G Y D N S P Y K M I E D L Y P N R F |
| BAS2924_1(41-115) | R E L S K R V T K Y L I E T I L K W N E E D I K Q K W N T P L I I K Y R L L G A L K H G Y D N S P Y K M I E D L Y P N R F |
| RBTH_06405_1(34-108) | N Q L A R R V T K Y L V T K I L N W N E E E I K Q N W N N K L I A K Y R L R G V L K H K Y N N S P Y A M I N D L Y P N Q F |
| pE33L466_0092_1(34-108) | N K M A R R V L T Y L L N S I L K W N K E D I R K K W N T K L L V K Y R L R G L L K H R Y E N S P F K A I N D L Y P N Q F |
| | |
| consensus/80% | + E K A L p s L + W h l E c c E K l s . . p L h p h a s . K W L . p . p L . u s h . h h W s s S P Y t M I N s L Y P s p a |

| Secondary structure | | |
|---|---|---|
| RBTH_06405_4(259-331) | I K S S F S G Y I N K F - - | 73 |
| pBMB165_3(175-247) | L K S S F R G Y I N K S - - | 73 |
| pE33L466_0092_4(259-328) | S K N M L K G Y M - - - - - | 70 |
| RBTH_06405_3(109-183) | K E W E F N K A P N K F W T | 75 |
| pBMB165_2(25-99) | K E W E F T K A P N K F W T | 75 |
| BA3147_2(109-183) | K E W E F S M T P N N F W T | 75 |
| Bant_01003795_1(25-99) | K E W E F S M T P N N F W T | 75 |
| BAS2924_2(116-190) | K E W E F S M T P N N F W T | 75 |
| BAS2924_3(191-265) | K E W E F S M T P N N F W T | 75 |
| pE33L466_0092_2(109-183) | K E W E F G M T P N N F W T | 75 |
| RBTH_06405_2(184-258) | K E W E F K M T P N G F W T | 75 |
| pBMB165_1(100-174) | K E W E F K M T P S G F W T | 75 |
| BAS2924_4(266-340) | K E W E Y K F T P T G F W T | 75 |
| Bant_01003795_2(100-174) | K E W E Y K F T P T G F W T | 75 |
| BA3147_3(184-258) | K E W E Y K F T P T G F W T | 75 |
| pE33L466_0092_3(184-258) | K E W E F R V T P V G Y W S | 75 |
| BA3147_1(34-108) | K E W E F G M A P L N F W T | 75 |
| BAS2924_1(41-115) | K E W E F G M A P L N F W T | 75 |
| RBTH_06405_1(34-108) | K E W E F R M T P L N F W T | 75 |
| pE33L466_0092_1(34-108) | K E W E F G M T P L N F W T | 75 |
| | |
| consensus/80% | K E W E F p h s P . t F W T |

FIGURE 11: BA3147 is homologous to protein GBAA3147 from *Bacillus anthracis* str. "*Ames* Ancestor."

| Secondary structure | E E E | HHHHHH | E E E | |
|---|---|---|---|---|
| BAS2851_1(20-78) | L E Y Q Q S R F Y V T R I P K D F L S I A R K R F S I P T D D Q I I A F L S C N L F G - - - S G K Y G V Y F T S S G L Y W K | 59 |
| BA3065_1(13-71) | L E Y Q Q S R F Y V T R I P K D F L S I A R K R F S I P T D D Q I I A F L S C N L F G - - - S G K Y G V Y F T S S G L Y W K | 59 |
| Bant_01003715_1(16-74) | L E Y Q Q S R F Y V T R I P K D F L S I A R K R F S I P T D D Q I I A F L S C N L F G - - - S G K Y G V Y F T S S G L Y W K | 59 |
| BcerKBAB4DRAFT_1832_1(14-72) | L E Y Q Q S R F Y V T R I P K D F L S V A K K R F S I P I D D R I F A F L S C N L F G - - - S G K Y G V Y F T S S G L Y W K | 59 |
| RBTH_02124_1(13-71) | L E F Q Q S R F Y V T R I P K D F L S I A Q K R F S I P T E D Q I V A F L S C N L L G - - - S G K Y G V Y F T S S G L Y W K | 59 |
| Bant_01003715_2(164-225) | L E P D N G L F V E T H I S D K K L K A I E V R F I I P I E E Q I I A F L D T S V L G N M G K G S D G V L I C Q S G I Y F R | 62 |
| BAS2851_2(168-229) | L E P D N G L F V E T H I S D K K L K A I E V R F I I P I E E Q I I A F L D T S V L G N M G K G S D G V L I C Q S G I Y F R | 62 |
| BA3065_2(161-222) | L E P D N G L F V E T H I S D K K L K A I E V R F I I P I E E Q I I A F L D T S V L G N M G K G S D G V L I C Q S G I Y F R | 62 |
| BcerKBAB4DRAFT_1832_2(162-223) | L E P D N G L F V D T H I S H K K L K E I G A K Y I I P K E E K I I A F L D T S V L G N L G K G S D G V L I C E P G I Y F R | 62 |
| | |
| consensus/80% | L E . p p u h F h . T + I s c c h L p h h p h R F . I P h - - p I l A F L s s s l h G . . . p G p . G V h h s p S G l Y a + |

FIGURE 12: BA3065 is homologous to protein GBAA3065 from *Bacillus anthracis* str. "*Ames* Ancestor."

| Secondary structure | E E E E | E E E | E E E E E | |
|---|---|---|---|---|
| BA0482(4-56) | I E I H T Q G G L K H K V Q T E V Y N A E A L N T K L N D N D L I T V L I G D F I I Q R I D V K R I I P L | 53 |
| BA0482(67-119) | V E V H T N A G K V I E I T T N D Y D P I Y L N E Q L N N N N T I T V V I G D Y I F S R I D V K Q V V P V | 53 |
| | |
| consensus/80% | l E l H T p u G h h h c l p T p s Y s s . h L N p p L N s N s h I T V l I G D a I h p R I D V K p l l P l |

FIGURE 13: BA0482 is homologous to proteins GBAA0482 from *Bacillus anthracis* str. "*Ames* Ancestor," BAS0458 from *Bacillus anthracis* str. Sterne, and Bant_01001108 from *Bacillus anthracis* str. A2012.
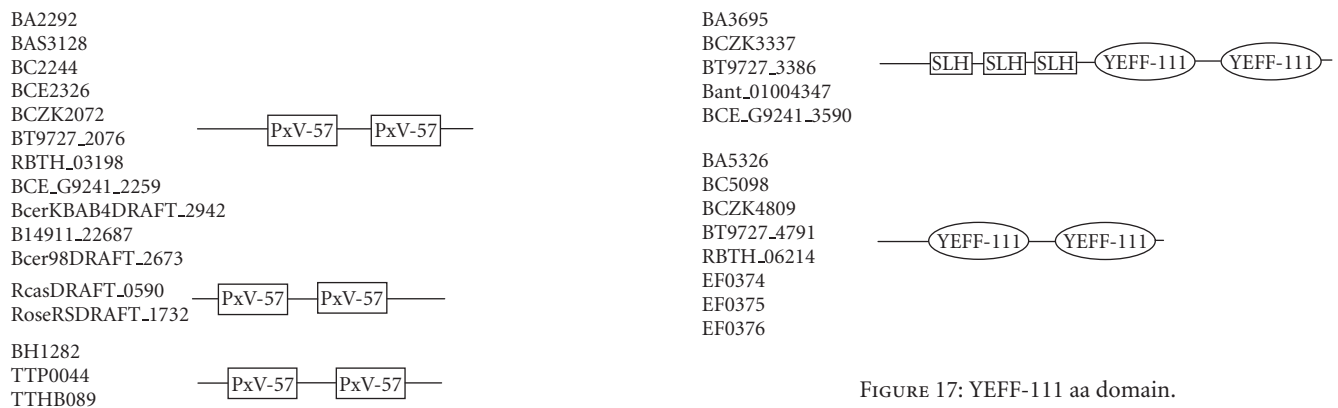
```
Secondary structure                        HHHH
BA4081(10-50)      S I GMY L S E L QK G T E S S R L L A E SMAK E I DGKMK I D L G P A G Q F    41
BA4081(172-212)    N I Q T L I NGMQ I G A L S L P Q V A Q TMG L D I K S N V Q V D L G E A G Q F  41

consensus/80%      s I t h h l s t h Q h G s . S . . . l A p o M u h - I c u p h p l D L G . A G Q F
```

(a)

```
Secondary structure                   E E E            HHH
BA4081(292-333)    G S K S G S E L G Q G I I S Q D G Y I K G S A L Q V V G S AHN A F S T I NG S P A  42
BA4081(334-375)    GN Q G G Q G F G S G I V N Q K G Y I R G S A L E A V T P AH T G F N T I NG T P Q  42

consensus/80%      G s p u G p t h G p G I l s Q c G Y I + G S A L p s V s s AH s u F s T I NG o P t
```

(b)

FIGURE 14: BA4081 is homologous to proteins GBAA4081 from *Bacillus anthracis* str. "*Ames* Ancestor," BAS3792 from *Bacillus anthracis* str. Sterne, and Bant_01004731 from *Bacillus anthracis* str. A2012.

BA2292
BAS3128
BC2244
BCE2326
BCZK2072
BT9727_2076
RBTH_03198
BCE_G9241_2259
BcerKBAB4DRAFT_2942
B14911_22687
Bcer98DRAFT_2673

RcasDRAFT_0590
RoseRSDRAFT_1732

BH1282
TTP0044
TTHB089

Chlo02001630
ExigDRAFT_0608
CaggDRAFT_2922

HaurDRAFT_2803

NT01CX_1619

SamaDRAFT_3539
rrnAC0576

Ava_3757

FIGURE 15: PxV-57 aa domain.

BA0881
BCZK0785
BT9727_0783
BCE_G9241_0886
Bcer98DRAFT_3031
B14911_04439
GK3171

CTC00525
DredDRAFT_0533
NT01CX_1557

FIGURE 16: FxF-122 aa domain.

BA3695
BCZK3337
BT9727_3386
Bant_01004347
BCE_G9241_3590

BA5326
BC5098
BCZK4809
BT9727_4791
RBTH_06214
EF0374
EF0375
EF0376

FIGURE 17: YEFF-111 aa domain.

BA1021
BAS0955
BC1029
BCZK0933
BT9727_0941
RBTH_03050
BCE_G9241_1042
BcerKBAB4DRAFT_3543
AmetDRAFT_1908
Bcer98DRAFT_1038
CAC3450
CPE0158
CbeiDRAFT_3331
DhafDRAFT_0725
CtheDRAFT_1311
CTC02189
CphyDRAFT_3436
ClosDRAFT_1658
CD1511
CPF_0149
CdifQ_02001573

BcerKBAB4DRAFT_0307

FIGURE 18: IMxxH-109 aa domain.

BA4716
BC449
BCE4587
BCZK4235
BT9727_4219     —( VxxT-103 )—( VxxT-103 )—
RBTH_05210
Bcer98DRAFT_3179
BcerKBAB4DRAFT_4089
ABF83609

BSU28380
BL00314
BH3070     —( VxxT-103 )—( VxxT-103 )—
B14911_06091

GK2667
OB2107     —( VxxT-103 )—( VxxT-103 )—

ABC2653     —( VxxT-103 )—( VxxT-103 )—

SwolDRAFT_2302
MothDRAFT_0979     —( VxxT-103 )—
Moth_0516
germ

CtheDRAFT_0840
GAA01614     —( VxxT-103 )—

AmetDRAFT_1640     —( VxxT-103 )—( VxxT-103 )—

FIGURE 19: VxxT-103 aa domain.

BA4310
BC4088
BCE4157
BCZK3845
BT9727_3829
BH0678
Bant_01004966
RBTH_02670     —( ExW-84 )—( ExW-84 )—
BCE_G9241_4093
BcerKBAB4DRAFT_2040
B14911_05359
GK0969
BAA83944
Bcer98DRAFT_3614
ExigDRAFT_0574

BSU30660
BL05305
BH0983
B14911_09907
BH1853     —( ExW-84 )—
ABC0230
ABC4088
OB3282
OB2488
ExigDRAFT_1796

FIGURE 20: ExW-84 aa domain.

BA2665
BC2674
BCE2700
BCZK2413
BT9727_2444     —[ NTGFIG-104 ]₂—
Bant_01003317
BCE_G9241_CNI_0263
BcerKBAB4DRAFT_0535
Bcer98DRAFT_0128

FIGURE 21: NTGFIG-104 aa domain.

copies. Further BLAST searches using sequence corresponding to the region (94–129) as query identified 9 hypothetical proteins comprising this repeat region from the organisms *B. anthracis*, *B. thuringiensis*, *B. thuringiensis serovar israelensis, Bacillus weihenstephanensis*, and *B. cereus* (see Table 1(h)). The length of proteins varied between 189 to 193-amino-acid residues, and also consists a SAP domain at the N-terminus, in addition to the novel repeat described here. A SAP domain consists of two $\alpha$-helices and is a DNA-binding motif that is involved in chromosomal organization [32]. Therefore, we believe that these repeats might also participate in a similar function. The multiple sequence alignment corresponding to this repeat identified NxGK sequence motif (Figure 9). The pairwise sequence identities between sequences corresponding to NxGK repeats varied between 36–97%. The secondary structure is predicted to comprise a $\alpha$-helix and the conserved sequence motif described above is also associated with $\alpha$-helix. The representative domain architecture corresponding to proteins comprising the NxGK repeats is shown in Figure 22.

### 3.9. 95-amino-acid-residue VYV domain

The 225-amino-acid-residue protein corresponding to the GENE_ID BA1701 and described as a hypothetical protein comprises a 95-amino-acid-residue region, as two copies in tandem. Further BLAST searches using sequence corresponding to the region (31–125) as query identified BAS1577 protein of *B. anthracis*, RBTH_03882 protein of *Bacillus*

*thuringiensis serovar israelensis*, and DSY3134 of *Desulfitobacterium hafniense* Y51 that are described as hypothetical proteins. The length of proteins varied between 227 to 1674-amino-acid residues (see Table 1(i)). In RBTH_03882, this region occurs ten times and in tandem. The multiple sequence alignment corresponding to this domain identified characteristic sequence motifs; GDxV, VYV (see Figure 10). For the sake of simplicity, we refer to this 95-amino-acid region as VYV domain. The pairwise sequence identities between sequences corresponding to VYV domains varied between 29–95%. The secondary structure corresponding to VYV domain is predicted to comprise five $\beta$-strands. The representative domain architecture corresponding to proteins comprising the VYV domains is shown in Figure 23.

### 3.10. 75-amino-acid-residue KEWE domain

The 262-amino-acid-residue protein corresponding to the GENE_ID BA3147 and described as a hypothetical protein comprises a 75-amino-acid-residue region as three copies in

BA3686
BC3626
BCE3645
BCZK3328
RBTH_03615
BT9727_3378
BCE_G9241_3579
BcerKBAB4DRAFT_0944
B14911_25780

Figure 22: NxGK-36 aa repeat.



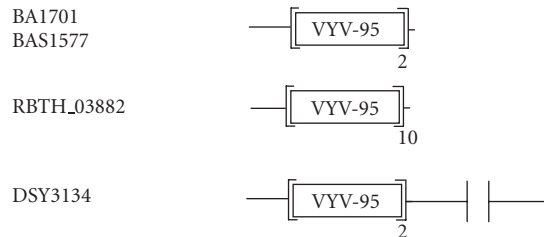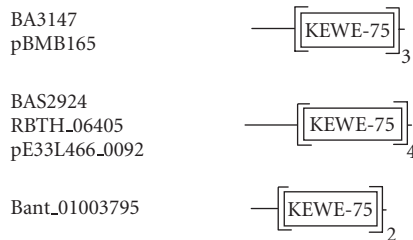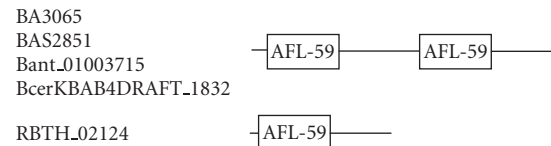BA1701
BAS1577

RBTH_03882

DSY3134

Figure 23: VYV-95 aa domain.



BA3065
BAS2851
Bant_01003715
BcerKBAB4DRAFT_1832

RBTH_02124

Figure 25: AFL-59 aa domain.



BA3147
pBMB165

BAS2924
RBTH_06405
pE33L466_0092

Bant_01003795

Figure 24: KEWE-75 aa domain.

tandem. Further BLAST searches using the sequence corresponding to the region (34–108) as query identified this domain in 6 proteins that are described as hypothetical proteins (see Table 1(j)). This domain may exist as 2, 3, or 4 copies in these proteins. The length of proteins identified varied between 178 to 344-amino-acid residues. The pairwise sequence identities between sequences corresponding to these regions varied between 22–69%. These domains are present in tandem and associated with SPY, MIN, LYP, KEWE, and FWT conserved sequence motifs as indicated in the multiple sequence alignment (see Figure 11). We refer to these as the KEWE domain, and this sequence motif occurs at the C-terminus of the domain. The secondary structure corresponding to KEWE domain is predicted to comprise three $\alpha$-helices as shown in Figure 11. The representative domain architecture corresponding to proteins comprising the KEWE domain is shown in Figure 24.

### 3.11. 59-amino-acid-residue AFL domain

The 290-amino-acid-residue protein corresponding to the GENE_ID BA3065 and described as hypothetical protein comprises a 59-amino-acid-residue region as two copies.

Further BLAST searches using sequence corresponding to the region (13–71) as query identified that this region occurs twice in the proteins with GENE_ID's: BAS2851 and Bant_01003715 of *B. anthracis* strains, the protein with GENE_ID: BcerKBAB4DRAFT_1832 of *Bacillus weihenstephanensis,* and once in the protein with GENE_ID: RBTH_02124 of *Bacillus thuringiensis serovar israelensis* (see Table 1(k)). The lengths of the proteins varied between 145 to 297-amino-acid residues and are described as hypothetical proteins. The multiple sequence alignment corresponding to this domain identified two characteristic sequence motifs: RFxI and AFL (see Figure 12). We refer to this as the AFL domain. The sequence identities shared between AFL domains varied between 38–91%. The secondary structure corresponding to the AFL domain is predicted to comprise of one $\alpha$-helix and two $\beta$-strands and the conserved sequence motif AFL is a part of the $\alpha$-helix. The representative domain architecture corresponding to protein comprising the AFL domain is shown in Figure 25.

### 3.12. 53-amino-acid-residue RIDVK repeat

The 159-amino-acid-residue protein corresponding to the GENE_ID BA0482 and described as a conserved domain protein comprises a 53-amino-acid region as two copies. BLAST did not identify this repeat in any other proteins; therefore this repeat is unique to *B. anthracis* str. *Ames*. The multiple sequence alignment corresponding to this repeat identified three characteristic sequence motifs: ITV, IGD, and RIDVK (Figure 13). We refer to this as the RIDVK repeat. The sequence identity shared between this RIDVK repeats in BA0482 is 45%. The secondary structure corresponding to the RIDVK repeat is predicted to comprise three $\beta$-strands. The representative domain architecture corresponding to protein comprising the RIDVK repeat is shown in Figure 26.

FIGURE 26: RIDVK-53 aa repeat.



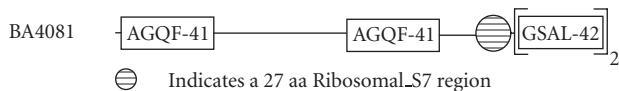⊖ Indicates a 27 aa Ribosomal_S7 region

FIGURE 27: AGQF-41 aa repeat; GSAL-42 aa repeat.

### 3.13. (a) 41-amino-acid-residue AGQF repeat and (b) 42-amino-acid-residue GSAL repeat

The protein corresponding to the GENE_ID BA4081 comprises 462-amino-acid residues and described as conserved domain protein contains two novel repeat types. The sequence length corresponding to repeat types are 41 and 42-amino-acid residues and are present as two copies in BA4081. BLAST searches identified these repeats to be specific to this protein alone.

(a) The sequence alignment corresponding to 41-amino-acid-residue repeat identified two characteristic sequence motifs: DLG and AGQF (Figure 14(a)). We refer to this as the AGQF repeat. The motif occurs at the C-terminal part of the repeat region. The sequence homology shared between this AGQF repeats is about 34%. The secondary structure corresponding to the AGQF repeat is predicted to comprise one $\alpha$-helix. The representative domain architecture corresponding to protein comprising the AGQF repeat is shown in Figure 27.

(b) The sequence alignment corresponding to the 42-amino-acid-residue tandem repeat identified three characteristic sequence motifs: GYI, GSAL, and TING (Figure 14(b)) and is a glycine-rich repeat. We refer to this as the GSAL repeat. The sequence homology shared between this GSAL repeats is 52%. The secondary structure corresponding to the GSAL repeat is predicted to comprise one $\alpha$-helix and one $\beta$-strand. The representative domain architecture corresponding to protein comprising the GSAL repeat is shown in Figure 27. This protein is associated with a 27-amino-acid-residue Ribosomal_S7 region that is sandwiched between the 41-amino-acid-residue AGQF repeat and the 42-amino-acid-residue GSAL repeat. These two repeats are specific to this protein alone and are therefore *B. anthracis* str. *Ames* specific.

From the analysis of the *B. anthracis* proteome, we observed that the novel repeats and domains are present in all the strains, such as *Ames*, *Ames* ancestor, Sterne, and A2012, that have been sequenced so far. This indicates that these strains of *B. anthracis* have diverged recently. We also observed that the domains PxV, FxF, YEFF, VxxT, ExW, and VYV are present in proteins from several bacterial organisms. The domains NTGFIG, KEWE, AFL, and the repeats NxGK are specific to bacillus. It is interesting to note that the domains VYV and AFL are present in all the *B. anthracis* species

while absent in *B. cereus* genomes. The repeats RIDVK, AGQF, and GSAL are also specifically present only in all the strains of *B. anthracis*. This analysis explains some differences in the closely related *B. anthracis* and *B. cereus* genomes. The identification of these novel domains and repeats in subsequently sequenced genomes will add value to their annotation.

## 4. CONCLUSIONS

A systematic analysis using computational tools identified four novel repeats and ten domains corresponding to the *B. anthracis* str. *Ames* proteome. Further database searches identified that some novel repeats and domains are also present in other bacterial genomes. The NxGK repeats are associated with SAP domain. The SAP domain is a DNA-binding motif that is involved in chromosomal organization. Therefore, we believe that these repeats also participate in similar function. The YEFF domain-containing proteins are associated with RGD motif and may be involved in cell adhesion. The identification of novel repeats and domains corresponding to *B. anthracis* proteome may be useful for annotation. From the presence of VYV and AFL domains in all the *B. anthracis* species and their absence in *B. cereus* genomes, we identified some differences in these two genomes that are otherwise closely related.

## REFERENCES

[1] R. Okinaka, K. Cloud, O. Hampton, et al., "Sequence, assembly and analysis of pX01 and pX02," *Journal of Applied Microbiology*, vol. 87, no. 2, pp. 261–262, 1999.

[2] T. C. Dixon, M. Meselson, J. Guillemin, and P. C. Hanna, "Anthrax," *New England Journal of Medicine*, vol. 341, no. 11, pp. 815–826, 1999.

[3] I. Uchida, T. Sekizaki, K. Hashimoto, and N. Terakado, "Association of the encapsulation of Bacillus anthracis with a 60 megadalton plasmid," *Journal of General Microbiology*, vol. 131, no. 2, pp. 363–367, 1985.

[4] I. Uchida, K. Hashimoto, and N. Terakado, "Virulence and immunogenicity in experimental animals of Bacillus anthracis strains harbouring or lacking 110 MDa and 60 MDa plasmids," *Journal of General Microbiology*, vol. 132, no. 2, pp. 557–559, 1986.

[5] S. H. Leppla, "Anthrax toxins," in *Bacterial Toxins and Virulence Factors in Disease*, J. Moss, B. Iglewski, M. Vaughn, and A. T. Tu, Eds., pp. 543–572, Marcel Dekker, New York, NY, USA, 1995.

[6] M. J. Pallen, A. C. Lam, M. Antonio, and K. Dunbar, "An embarrassment of sortases - a richness of substrates?" *Trends in Microbiology*, vol. 9, no. 3, pp. 97–101, 2001.

[7] J. M. Patti, B. L. Allen, M. J. McGavin, and M. Hook, "MSCRAMM-mediacted adherence of microorganisms to host tissues," *Annual Review of Microbiology*, vol. 48, pp. 585–617, 1994.

[8] W. W. Navarre and O. Schneewind, "Surface proteins of gram-positive bacteria and mechanisms of their targeting to the cell wall envelope," *Microbiology and Molecular Biology Reviews*, vol. 63, no. 1, pp. 174–229, 1999.

[9] D. M. Guttmann and D. J. Ellar, "Phenotypic and genotypic comparisons of 23 strains from the Bacillus cereus complex for a selection of known and putative *B. thuringiensis* virulence factors," *FEMS Microbiology Letters*, vol. 188, no. 1, pp. 7–13, 2000.

[10] A. Lupas, H. Engelhardt, J. Peters, U. Santarius, S. Volker, and W. Baumeister, "Domain structure of the *Acetogenium kivui* surface layer revealed by electron crystallography and sequence analysis," *Journal of Bacteriology*, vol. 176, no. 5, pp. 1224–1233, 1994.

[11] A. Lupas, "A circular permutation event in the evolution of the SLH domain?" *Molecular Microbiology*, vol. 20, no. 4, pp. 897–898, 1996.

[12] T. D. Read, S. N. Peterson, N. Tourasse, et al., "The genome sequence of *Bacillus anthracis Ames* and comparison to closely related bacteria," *Nature*, vol. 423, no. 6935, pp. 81–86, 2003.

[13] I. Letunic, R. R. Copley, S. Schmidt, et al., "SMART 4.0: towards genomic data integration," *Nucleic Acids Research*, vol. 32, pp. D142–D144, 2004.

[14] N. J. Mulder, R. Apweiler, T. K. Attwood, et al., "The InterPro database, 2003 brings increased coverage and new features," *Nucleic Acids Research*, vol. 31, no. 1, pp. 315–318, 2003.

[15] A. Bateman, E. Birney, L. Cerruti, et al., "The Pfam protein families database," *Nucleic Acids Research*, vol. 30, no. 1, pp. 276–280, 2002.

[16] X. Yu, C. C. S. Chini, M. He, G. Mer, and J. Chen, "The BRCT domain is a phospho-protein binding domain," *Science*, vol. 302, no. 5645, pp. 639–642, 2003.

[17] B. Kobe and J. Deisenhofer, "The leucine-rich repeat: a versatile binding motif," *Trends in Biochemical Sciences*, vol. 19, no. 10, pp. 415–421, 1994.

[18] M. A. Andrade, C. Perez-Iratxeta, and C. P. Ponting, "Protein repeats: structures, functions, and evolution," *Journal of Structural Biology*, vol. 134, no. 2-3, pp. 117–131, 2001.

[19] M. A. Andrade, F. D. Ciccarelli, C. Perez-Iratxeta, and P. Bork, "NEAT: a domain duplicated in genes near the components of a putative $Fe^{3+}$ siderophore transporter from Gram-positive pathogenic bacteria," *Genome Biology*, vol. 3, no. 9, pp. 0047.1–0047.5, 2002.

[20] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, "Basic local alignment search tool," *Journal of Molecular Biology*, vol. 215, no. 3, pp. 403–410, 1990.

[21] A. Heger and L. Holm, "Rapid automatic detection and alignment of repeats in protein sequences," *Proteins: Structure, Function and Genetics*, vol. 41, no. 2, pp. 224–237, 2000.

[22] M. A. Andrade, C. P. Ponting, T. J. Gibson, and P. Bork, "Homology-based method for identification of protein repeats using statistical significance estimates," *Journal of Molecular Biology*, vol. 298, no. 3, pp. 521–537, 2000.

[23] R. Mott, "Accurate formula for P-values of gapped local sequence and profile alignments," *Journal of Molecular Biology*, vol. 300, no. 3, pp. 649–659, 2000.

[24] R. Szklarczyk and J. Heringa, "Tracking repeats using significance and transitivity," *Bioinformatics*, vol. 20, supplement 1, pp. i311–i317, 2004.

[25] S. F. Altschul, T. L. Madden, A. A. Schäffer, et al., "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs," *Nucleic Acids Research*, vol. 25, no. 17, pp. 3389–3402, 1997.

[26] J. D. Thompson, D. G. Higgins, and T. J. Gibson, "CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice," *Nucleic Acids Research*, vol. 22, no. 22, pp. 4673–4680, 1994.

[27] B. Rost, C. Sander, and R. Schneider, "PHD—an automatic mail server for protein secondary structure prediction," *Computer Applications in the Biosciences*, vol. 10, no. 1, pp. 53–60, 1994.

[28] L. Falquet, M. Pagni, P. Bucher, et al., "The PROSITE database, its status in 2002," *Nucleic Acids Research*, vol. 30, no. 1, pp. 235–238, 2002.

[29] C. P. Ponting, R. Mott, P. Bork, and R. R. Copley, "Novel protein domains and repeats in *Drosophila melanogaster*: insights into structure, function, and evolution," *Genome Research*, vol. 11, no. 12, pp. 1996–2008, 2001.

[30] S. M. Akula, N. P. Pramod, F.-Z. Wang, and B. Chandran, "Integrin $\alpha3\beta1$ (CD 49c/29) is a cellular receptor for Kaposi's sarcoma-associated herpesvirus (KSHV/HHV-8) entry into the target cells," *Cell*, vol. 108, no. 3, pp. 407–419, 2002.

[31] S. E. D'Souza, M. H. Ginsberg, and E. F. Plow, "Arginyl-glycyl-aspartic acid (RGD): a cell adhesion motif," *Trends in Biochemical Sciences*, vol. 16, no. 7, pp. 246–250, 1991.

[32] L. Aravind and E. V. Koonin, "SAP—a putative DNA-binding motif involved in chromosomal organization," *Trends in Biochemical Sciences*, vol. 25, no. 3, pp. 112–114, 2000.