# Identification and characterization of a cyanobacterial DnaX intein

Xiang-Qin Liu*, Zhuma Hu

*Biochemistry Department, Dalhousie University, Halifax, Nova Scotia B3H 4H7, Canada*

**Abstract** A new intein is identified and characterized in the DnaX protein of *Synechocystis* sp. PCC6803. This cyanobacterial DnaX protein is a homologue of the intein-less 71-kDa *tau*-subunit of *Escherichia coli* DNA polymerase III and is related to eukaryotic DNA replication factor C (RFC). The 430-residue DnaX intein contains several putative intein sequence motifs and undergoes protein splicing when produced in *E. coli* cells. Its position in the DnaX protein is close to, but different from, positions of three inteins present in a DnaX-related RFC protein of *Methanococcus jannaschii*.

© 1997 Federation of European Biochemical Societies.

*Key words:* Intein; Protein splicing; DnaX; DNA polymerase subunit

## 1. Introduction

An intein is a protein sequence embedded inframe within a precursor protein sequence that is spliced out during maturation [1]. This maturation process is termed protein splicing and involves the precise excision of the intein sequence and the formation of a normal peptide bond joining the external sequences (exteins). Protein splicing is therefore considered the protein equivalent of RNA splicing and adds another layer of complexity to the Central Dogma of molecular biology. In-teins have been found in eukaryotes [2,3], eubacteria [4–8], and archaebacteria [1,9–11], suggesting a wide distribution.

There is little overall sequence similarity among inteins, unless comparison is made between inteins found at same position in homologous proteins of different organisms. Nevertheless, a number of short sequence motifs have been recognized that show a significant degree of conservation among inteins [7]. Protein splicing has been demonstrated experimentally for a small fraction of the reported inteins [2–5,9,12]. The chemistry of protein splicing has been discussed in several models and in some cases tested experimentally [10,12–15]. It generally involves N→S acyl shift (or N→O shift) at the splice sites [16], formation of a branched inter-mediate [17], and cyclization of an invariant Asn residue at the C-terminus of intein to succinimide [18], leading to exci-sion of the intein and ligation of the exteins.

Among the more than 30 intein or intein-like sequences reported, over 20 of them are in archaebacteria (archaea) [1,9–11], eight are in mycobacteria [4–7], with the remaining few in other organisms [2,3,7,8]. One intein coding sequence was recognized previously in a cyanobacterial DNA helicase gene *dnaB* [8], although it has not been characterized for pro-tein splicing. Here we report the identification and character-ization of a new cyanobacterial intein encoded in the *dnaX* gene of *Synechocystis* sp. strain PCC6803. The *dnaX* gene encodes a homologue of the *E. coli* intein-less 71-kDa *tau*-

subunit of DNA polymerase III. The cyanobacterial DnaX intein is compared with related inteins including those present in a DnaX-related DNA replication factor C protein. Protein splicing activity of the cyanobacterial DnaX intein is demon-strated in *E. coli* cells. This new intein has been registered with Dr. Perler at New England Biolabs where a complete list of known inteins is maintained (see ref. [1]).

## 2. Materials and methods

### 2.1. DNA sequence analysis and cloning

GenBank searches were performed using the BLAST search pro-gram [19]. Protein sequence alignment was carried out using the Clus-tal W program [20]. The 3351-bp *dnaX* coding sequence was prepared from *Synechocystis* sp. strain PCC6803 by PCR-amplification, using the thermostable DNA polymerase Pfu (Stratagene) and a pair of oligonucleotide primers: 5′-ATCGTGCCATGGCCTACGAACCTC-TG-3′ and 5′-ATCGAAAAGCAGTGAGTCCCATTAG-3′. Re-combinant plasmid pTX-1 was constructed by cloning a 1468-bp *Cla*I–*Nco*I DNA fragment of the *dnaX* gene into the expression plas-mid vector pET-32 (Novagen) between its *Bst*BI and *Nco*I sites. Re-combinant plasmid pTX-2 was constructed by cloning a 1695-bp DNA fragment of the *dnaX* gene into pET-32 between its *Bgl*II and *Bam*HI sites. pET-32 encodes a gene expression cassette: T7 pro-moter/thioredoxin-coding sequence/poly-histidine and S tag coding sequence/multiple cloning sites. In each case, the *dnaX* coding se-quence was fused inframe with the vector's coding sequence, placing the fusion gene directly behind an IPTG-inducible T7 promoter.

### 2.2. Protein production and splicing in E. coli cells

Each recombinant plasmid was introduced into *E. coli* strain DE3 that harbors an IPTG-inducible T7 RNA polymerase gene. *E. coli* cells transformed with individual recombinant plasmids of interest were grown in liquid LB medium at 37°C to late log phase ($A_{600} = 0.5$). IPTG was added at this time to a final concentration of 1 mM to induce production of the recombinant proteins, and the induction was continued for 3 h at 25°C. Before SDS-polyacrylamide gel electrophoresis, cells were lysed in SDS-containing gel loading buffer in a boiling water bath. For isolating proteins containing a poly-histidine tag (e.g. protein product of pTX-1), cells were lysed in a denaturing buffer (50 mM $NaH_2PO_4$, 10 mM Tris-HCl, 8 M urea, pH 8.0). These proteins were selectively precipitated by using the TALON metal affinity resin (from Clontech) which binds the poly-histidine peptide sequence. To selectively detect proteins containing the S tag, Western blotting was carried out by using an S protein (Novagen) that specifically recognizes the S tag sequence.

## 3. Results and discussion

### 3.1. Identification and analysis of the DnaX intein sequence

In a routine BLAST search [19] of the GenBank for homo-logues of a DnaB intein, we noticed putative intein sequence motifs encoded in the *dnaX* gene of the cyanobacterium *Syn-echocystis* sp. strain PCC6803 (*S.* sp.). The *dnaX* gene se-quence (GenBank Accession No. D90907) was reported re-cently by others as a part of the complete genome sequence of this organism [21]. Further analysis of this gene shows that the predicted *S.* sp. DnaX protein is much larger than its *E. coli* homologue, and this is almost entirely due to the presence of a large intervening sequence (430 amino acid residues long)

*Corresponding author. Fax: (902) 494-1355. E-mail: pxqliu@is.dal.ca

**A**

N-extein     **intein**     C-extein

Ssp

129    430      557

Eco

127      516

**B**

```
Ssp dnaX   MAYEPLHHKYRPQTFADLVGQTAIAATLSNAIEQERIVPA   40
           I:I: I :I IIIIIII:III  : ::I:I :  II I
Eco dnaX   MSYQVLARKWRPQTFADVVGQEHVLTALANGLSLGRIHHA   40

Ssp dnaX   YLFTGPRGTGKTSSARILAKSLNCIAGDRPTATPCGQCAT   80
           III:I II IIII II:III III :I  IIIIII I :
Eco dnaX   YLFSGTRGVGKTSIARLLAKGLNCETG--ITATPCGVCDN   78

Ssp dnaX   CRAITNGSALDVIEIDAASNTGVDNIREIIERAQFAPVQC  120
           II I :I  :I:IIIIIII I :I  I::::  I:II :
Eco dnaX   CREIEQGRFVDLIEIDAASRTKVEDTRDLLDNVQYAPARG  118
```

**Ssp dnaX intein**
**430 a.a.**

▼

```
Ssp dnaX   RYKVYVIDECHMLSTAAFNALLKTLEEPPERVVFVLATTD  590
           I:III:III IIII  :IIIIIIIIIIII:I I:IIIII
Eco dnaX   RFKVYLIDEVHMLSRHSFNALLKTLEEPPEHVKFLLATTD  158

Ssp dnaX   PQRVLPTIISRCQRFDYRRIPLQAMVDHLRYIAGRENINI  630
           II::  II:III :I  : : ::: :I  I : I:I
Eco dnaX   PQKLPVTILSRCLQFHLKALDVEQIRHQLEHILNEEHIAH  198

Ssp dnaX   DQPALTLVAQIANGGLRDAESLLDQ(655)......
           : II I:I: I I IIII II II
Eco dnaX   EPRALQLLARAAEGSLRDALSLTDQ(223)......
```

**C**

```
Ssp dnaX   CLTGDSQVLTRNG----LMSIDNPQIKGREVLSYNETLQQ   36
           I:: ::I   I  : I I:: I IIII  I IIII
Mle recA   CMNYSTRVTLADGSTEKIGKIVNNKMDVR-VLSYDPVTDR   39
                  motif A

Ssp dnaX   WEYKKVLRWLDRG-EKQTLSIKTK------NSTVRCTANH   69
           :II: I  I  I I          I  :I II
Mle recA   IVPRKVVNWFNNGPAEQFLQFTVEKSGSNGKSQFAATPNH   79
                                            motif B

Ssp dnaX   LIRTEQGWTRAENITPGMKILSPASVDVDNLSQSTALTAS  109
           IIII III I I: I ::I:  I:  II
Mle recA   LIRTPGGWTEAGNLIAGDRVLA---VEPHMLSDQ-QFQVV  115

Ssp dnaX   LGGLSGAINY....(285 a.a. residues)....GQ  406
           II I I I                             :
Mle recA   LGSLMGDGNL....(214 a.a. residues)....TR  341
             motif C

Ssp dnaX   VEKVYDLEVEDNHNFVANGLLVHN   430
           :I:III III:  I::III
Mle recA   SMNRFDIEVEGNHNYFVDGVMVHN   365
             motif F     motif G
```
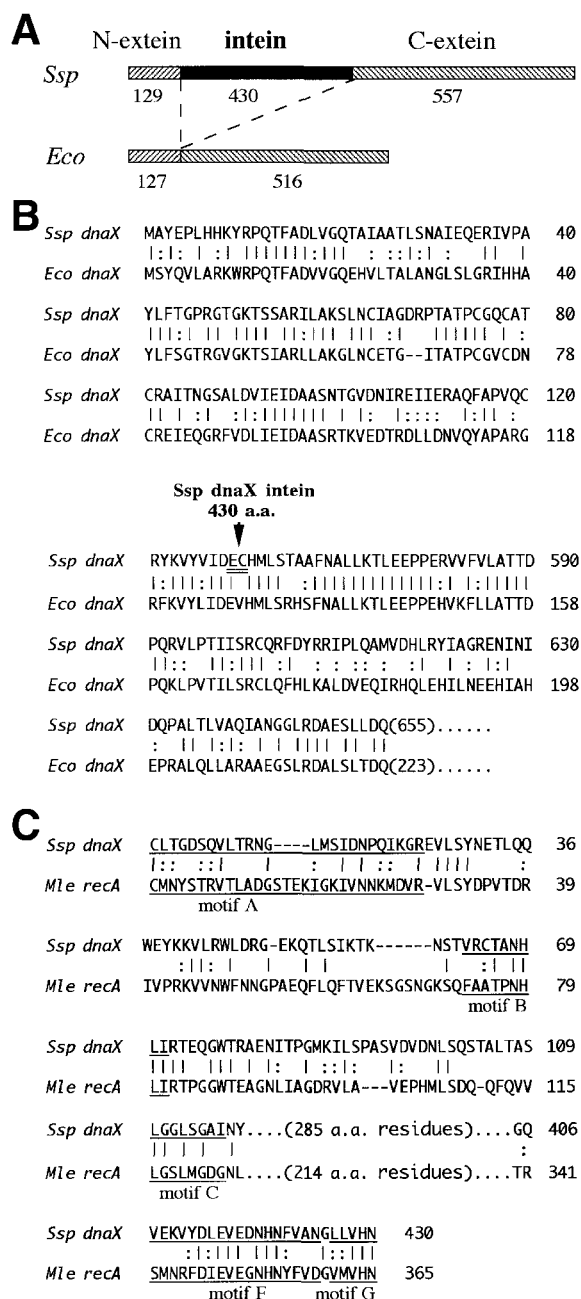
Fig. 1. Analysis of the DnaX intein sequence. A: Schematic illustration and comparison of DnaX protein between *Synechocystis* sp. PCC6803 (*Ssp*) and *E. coli* (*Eco*). Regions corresponding to the N- and C-exteins (hatched boxes) and intein (solid box) are marked, with the number of amino acid residues shown for each region. B: Comparison of DnaX sequence between *Synechocystis* sp. PCC6803 (*Ssp dnaX*) and *E. coli* (*Eco dnaX*). Only the N-terminal regions of the two sequences are shown. Position of intein in the *Ssp dnaX* sequence is marked by an arrow head, the intein size is shown, and the two residues flanking the intein are double underlined. C: Sequence comparison between *Synechocystis* sp. PCC6803 DnaX intein (*Ssp dnaX*) and *Mycobacterium lepre* recA intein (*Mle recA*). Only the terminal regions of the two sequences are compared and shown, with the number of uncompared residues shown in parenthesis. Putative intein motifs (A, B, C, F, G) are underlined. Symbols: – represents gaps introduced to optimize the alignment; I and : mark positions of identical and similar amino acids, respectively. GenBank Accession Numbers for *Synechocystis* sp. PCC6803 DnaX, *E. coli* DnaX, and *Mycobacterium lepre* recA are D90907, D90754, and X73822, respectively.

in *S.* sp. DnaX (Fig. 1A). This intervening sequence is subsequently called an intein, with its flanking sequences termed N- and C-exteins (Fig. 1A). The extein sequences can be aligned with sequences of other known DnaX proteins and, for example, shows 53% sequence identity to the N-terminal 223-residue sequence of *E. coli* DnaX protein (Fig. 1B).

The *S.* sp. DnaX intein, positioned between residue 129 and residue 559 of the whole sequence, interrupts a sequence region that is highly conserved among DnaX proteins, so that the intein boundaries are easily defined (Fig. 1B). The intein boundaries are also based on other intein-defining features including a nucleophilic residue (Cys) at its N-terminus, a His–Asn sequence at its C-terminus, and another nucleophilic residue (Cys) at the beginning of C-extein (Fig. 1B, 1C). These four residues and their positions are highly conserved among inteins and are known to be critical in the chemistry of protein splicing [16–18].

Comparison between the *S.* sp. DnaX intein and other known inteins revealed only low levels of sequence similarity in the N- and C-terminal regions. The highest similarity was found with the *Mycobacterium lepre* recA intein with 32% sequence identity over a 156-residue aligned region (Fig. 1C), while the remaining sequences of the two inteins show little or no similarity. It is therefore unlikely that the *S.* sp. DnaX intein shared a recent origin with any of the other known inteins. Nevertheless, the *S.* sp. DnaX intein contains at least five putative intein motifs (motifs A, B, C, F, and G in Fig. 1C) out of the seven motifs conserved among other inteins [7]. This may suggest similarity between the *S.* sp. DnaX intein and other inteins in terms of protein splicing mechanism and evolutionary origin. In addition, the putative intein motif C is similar to the LAGLI-DADG motif found in many intein- and intron-encoded endonucleases, suggesting that this DnaX intein has or once had endonuclease activity.

### 3.2. Demonstration of protein splicing with the DnaX intein

The *S.* sp. DnaX intein was tested for protein splicing in *E. coli* cells. Two recombinant plasmids (pTX-1 and pTX-2) were constructed to encode fusion proteins consisting of *S.* sp. DnaX sequence and terminal tag sequences (Fig. 2A). Each fusion protein consists of the complete intein sequence flanked by various amount of extein sequences and tag sequences. Each recombinant plasmid was introduced into *E. coli* cells to produce the corresponding fusion protein and to observe possible protein splicing products (Fig. 2B). Cells containing plasmid pTX-1 produced three proteins with apparent sizes of 68 kDa, 51 kDa, and 24 kDa, which corresponded closely to sizes predicted for a precursor protein (72 kDa), an excised intein (48 kDa), and a spliced protein (24 kDa), respectively. Cells containing plasmid pTX-2 also produced three proteins having apparent sizes of 76 kDa, 51 kDa, and 35 kDa, respectively. Again, these sizes corresponded well with the predicted sizes of a precursor protein (82 kDa), an excised intein (48 kDa), and a spliced protein (34 kDa), respectively. In addition to size, the precursor protein and the spliced protein were further identified by their selective binding to metal affinity resin (property of the poly-histidine tag) and to the S protein (property of the S tag) as predicted. The intein band was identified by its size as well as by the fact that its size was not affected by differing extein sequences in pTX-1 and pTX-2.

The above results clearly demonstrate that the *S.* sp. DnaX intein is capable of protein splicing when produced in *E. coli*
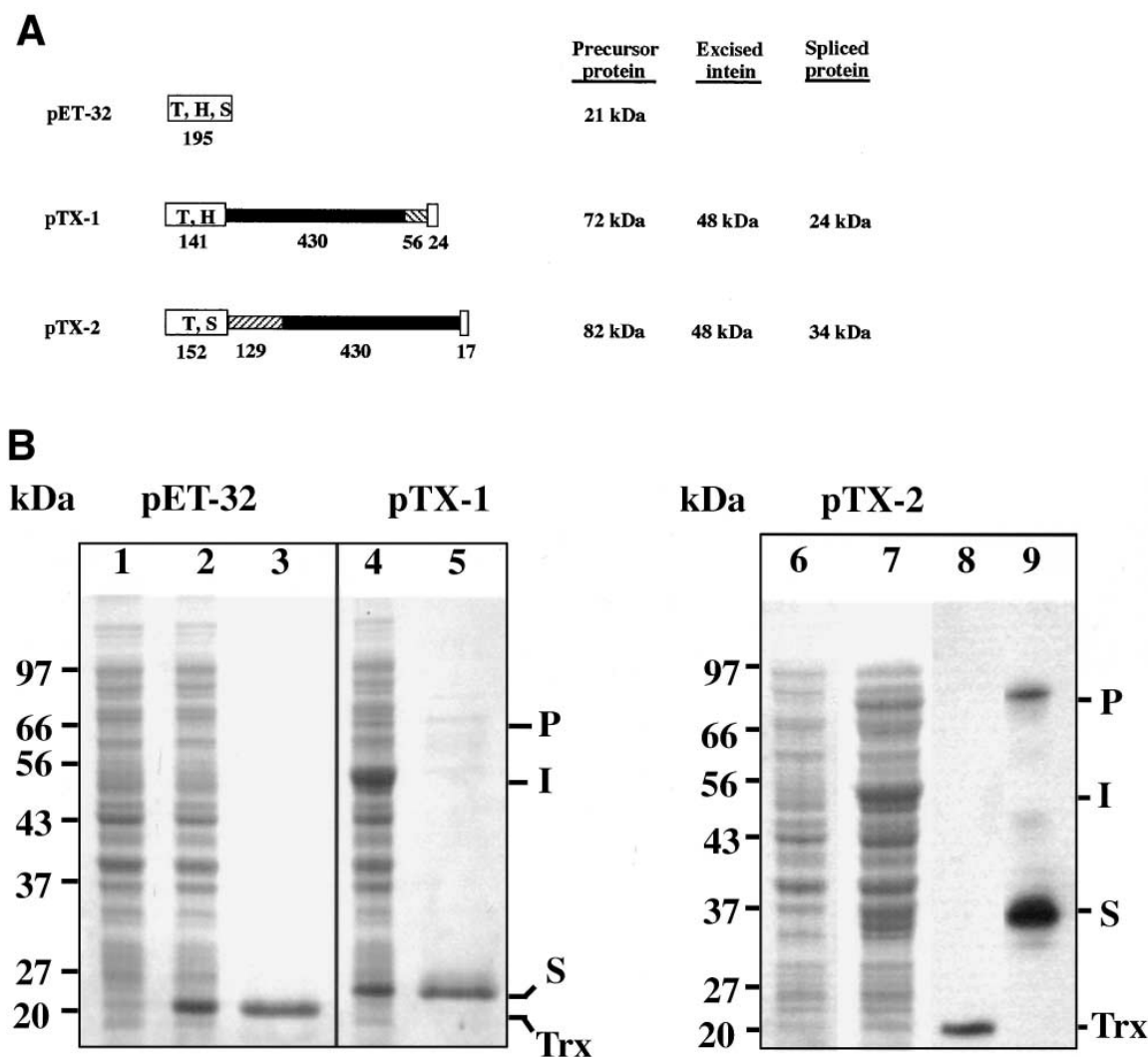
Fig. 2. Protein splicing of the DnaX intein. The *Synechocystis* sp. PCC6803 *dnaX* coding sequence was cloned into the expression plasmid vector pET-32 to produce recombinant fusion proteins and to observe protein splicing in *E. coli* cells. A: Schematic illustration of fusion protein construct. DnaX intein (solid box) and extein (hatched box) sequences are fused with vector-encoded sequences (open boxes) that include a thioredoxin protein (T), a poly-histidine tag (H) and an S tag (S) sequences, with the number of residues shown for each sequence domain. For each construct (pTX-1, pTX-2, and pET-32 as a control), calculated molecular weights are shown for the predicted precursor protein, the excised intein and the spliced protein. B: Production and protein splicing of recombinant DnaX fusion proteins in *E. coli*. Cells containing the individual recombinant constructs (specified above each box) were induced by IPTG to produce the corresponding recombinant proteins. Proteins were separated by SDS-polyacrylamide gel electrophoresis and visualized either by Coomassie Blue staining (lanes 1–7) or after Western blotting (lanes 8 and 9). Lanes 1 and 6, before induction by IPTG. Lanes 2, 4, and 7, after induction by IPTG. Lane 3 and lane 5, proteins isolated from lane 2 and lane 4, respectively, using metal affinity resin that recognizes the poly-histidine tag. Lane 8 and lane 9 correspond to lane 2 and lane 7, respectively, after Western blotting using the S protein that recognizes the S tag. Letters P, I, S, and Trx mark positions of precursor protein, excised intein, spliced protein, and thioredoxin protein, respectively.

cells, which is consistent with its intein-like sequence features. The results also show that complete extein sequences are not required for the protein splicing, because pTX-1 encodes only 2 residues of the N-extein, while pTX-2 encodes only 6 residues of the C-extein. It has been shown before that certain intein sequences are sufficient for protein splicing when placed within a foreign or target protein immediately before a nucleophilic residue, suggesting that all *cis* information required for protein splicing is contained within the intein [4,12,22]. We were unable to express the DnaX protein with the complete extein sequences, presumably due to toxicity of the protein to the host *E. coli* cells.

*3.3. Evolutionary considerations of the S. sp. DnaX intein*

This is the first time an intein has been found in a DnaX

protein. However, a search of the GenBank revealed significant sequence similarity between *S.* sp. DnaX and an intein-containing replication factor C (RFC) protein (Fig. 3). The *Methanococcus jannaschii* RFC 37 kDa protein, a homologue of RFC protein of eukaryotic DNA polymerase, contains three intein sequences [11]. In Fig. 3, the N-terminal sequence (extein only) of *S.* sp. DnaX is aligned with that of *M. jannaschii* RFC 37 kDa protein to compare their sequences and intein positions. The two sequences show 37% identity and 60% similarity over the 154 aligned positions, although the similarity is considerably lower over the rest of their sequences. The four inteins of these two proteins are not specifically related in sequence, suggesting independent origins. The *S.* sp. DnaX intein is located at a position different from, but close

```
                  (box II)                              (box III)
Ssp DnaX   MAYEPLHHKYRPQTFADLVGQTAIAATLSNAIEQERIVPAYLFTGPRGTGKTSSARILAK   60
           :  |  ||||:|  |:||| |   |   :|: : :| ||:|| | |||::| ||:
Mja RFC    MVIIMEKPWVEKYRPKTLDDIVGQDEIVKRLKKYVEK-KSMPHLLFSGPPGVGKTTAALCLAR   610

                                               Mja RFC-1 intein▲
                                                      548 a.a.

                        (box IV)
Ssp DnaX   SLNCIAGDRPTATPCGQCATCRAITNGSALDVIEIDAASNTGVDNIR-EIIERAQFAPVQCR-   121
           |                               :|: |: |:| || : :|: |:
Mja RFC    DLFGENWR--------------------DNFLELNASDERGIDVIRTKVKDFARTKPIGDVP   1089

                                         ▲Mja RFC-2 intein
                                          437 a.a.

           Ssp dnaX intein
              430 a.a.
               ▼
                                                    (box VII)
Ssp DnaX   YKVYVIDECHMLSTAAFNALLKTLEEPPERVVFVLATTDPQRVLPTIISRCQRF(605)....
           :|:  :||:  |::  | ||| :|:|    :  |:|::: | :::|  |  ||| |
Mja RFC    FKIIFLDESDALTADAQNALRRTMEKYSDVCRFILSCNYPSKIIPPIQSRCAVF(1686)...

           (box V)        (box VI)        ▲Mja RFC-3 intein
                                          543 a.a.
```
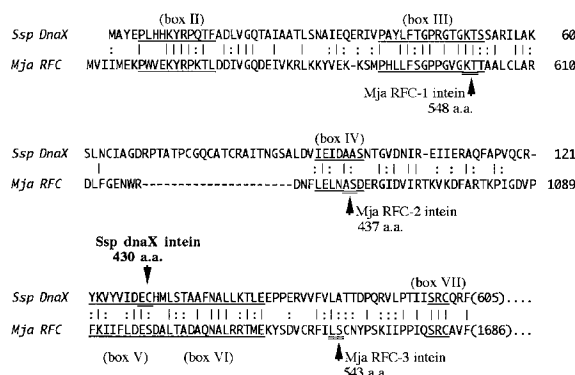
Fig. 3. Comparison between DnaX and RFC. The N-terminal sequence of *Synechocystis* sp. PCC6803 DnaX (*Ssp DnaX*) is aligned with that of *Methanococcus jannaschii* RFC 37 kDa protein (*Mja RFC*, Accession No. 1592072). The position and size of each intein are shown, with the two residues flanking each intein double underlined. Conserved sequence boxes II to VII, as defined previously in [23], are underlined. Symbols: − represents gaps introduced to optimize the alignment; ı and : mark positions of identical and similar amino acids, respectively.

to, the three intein positions of *M. jannaschii* RFC 37 kDa protein. The four inteins are clustered in a small sequence region close the N-termini of their respective proteins. This region of the sequence corresponds to an ATP/GTP-binding region that contains five small sequence boxes (boxes II to VI) conserved among a family of RFC-related proteins [23]. Interestingly, the *S.* sp. DnaX intein is positioned within a putative DEAD-box (box V, as defined in [23]), and each of the three *M. jannaschii* inteins is also positioned either within or adjacent to one of the conserved sequence boxes. Inteins positioned in a conserved region of a protein may be better maintained, because damage or inactivation of the intein would be less tolerated. Such an intein may also spread more easily through the intein homing process in which the intein coding sequence is copied into a new loci containing a suitable homing site [24–27], because a functional homing site would be more likely preserved in homologous proteins of different organisms if the intein is positioned in a conserved region of a protein.

*Synechocystis* sp. PCC6803 is the only eubacterium outside mycobacteria where intein is reported. The DnaX intein is one of only two inteins found in *Synechocystis* sp. PCC6803, with the other one encoded in the DNA helicase gene *dnaB* [8]. It is noted that the *E. coli* homologues of DnaX and DnaB are known to interact in DNA replication [28]. The two *S.* sp. inteins are not specifically related in sequence, which may suggest independent origins. The *E. coli* homologue of DnaX is the 71-kDa *tau*-subunit of the DNA polymerase III complex and has DNA-dependent ATPase and DNA annealing activities [29]. The finding of an intein-containing DnaX protein further highlights the interesting but unexplained observation that known inteins reside predominantly in proteins of DNA/RNA metabolism, which have previously included a DNA recombinase [4,5], a DNA polymerase [1,9,11], a DNA gyrase subunit [6], DNA helicases [8,11,30], a replication factor C protein [11], a reverse gyrase [11], RNA polymerase subunits [11], and a transcription initiation factor [11]. Less than 10 of the over 30 known inteins reside in other proteins [2,3,11].

## References

[1] F.B. Perler, E.O. Davis, G.E. Dean, F.S. Gimble, W.E. Jack, N. Neff, C.J. Noren, J. Thorner, M. Belfort, Nucl. Acids Res. 22 (1994) 1125–1127.
[2] P.M. Kane, C.T. Yamashiro, D.F. Wolczyk, N. Neff, M. Goebl, T.H. Stevens, Science 250 (1990) 651–657.
[3] H.H. Gu, J. Xu, M. Gallagher, G.E. Dean, J. Biol. Chem. 268 (1993) 7372–7381.
[4] E.O. Davis, P.J. Jenner, P.C. Brooks, M.J. Colston, S.G. Sedgwick, Cell 71 (1992) 201–210.
[5] E.O. Davis, H.S. Thangaraj, P.C. Brooks, M.J. Colston, EMBO J. 13 (1994) 699–703.
[6] H. Fsihi, V. Vincent, S.T. Cole, Proc. Natl. Acad. Sci. USA 93 (1996) 3410–3415.
[7] S. Pietrokovski, Protein Sci. 3 (1994) 2340–2350.
[8] S. Pietrokovski, Trends Genet. 12 (1996) 287–288.
[9] F.B. Perler, D.G. Comb, W.E. Jack, L.S. Moran, B. Qiang, R.B. Kucera, J. Benner, B.E. Slatko, D.O. Nwankwo, S.K. Hempstead, C.K.S. Carlow, H. Jannasch, Proc. Natl. Acad. Sci. USA 89 (1992) 5577–5581.
[10] R.A. Hodges, F.B. Perler, C.J. Noren, W.E. Jack, Nucl. Acids Res. 20 (1992) 6153–6157.
[11] Bult, C.J., White, W., Olsen, G.J., Zhouk, L., Fleischmann, R.D., Sutton, G.G., Blake, J.A., FitzGerald, L.M., Clayton, R.A., Gocayne, J.D., Kerlavage, A.R., Dougherty, B.A., Tomb, J.-F., Adams, M.D., Reich, C.I., Overbeek, R., Kirkness, E.F., Weinstock, K.G., Merrick, J.M., Glodek, A., Scott, J.L., Geoghagen, N.S.M., Weidman, J.F., Fuhrmann, J.L., Nguyen, D., Utterback, T.R., Kelley, J.M., Peterson, J.D., Sadow, P.W., Hanna, M.C., Cotton, M.D., Roberts, K.M., Hurst, M.A., Kaine, B.P., Borodovsky, M., Klenk, H.-P., Fraser, C.M., Smith, H.O., Woese, C.R. and Venter, J.C. (1996) Science 273, 1058-1073.
[12] M.Q. Xu, M.W. Southworth, F.B. Mersha, L.J. Hornstra, F.B. Perler, Cell 75 (1993) 1371–1377.
[13] C.J. Wallace, Protein Sci. 2 (1993) 697–705.
[14] N.D. Clarke, Proc. Natl. Acad. Sci. USA 91 (1994) 11084–11088.
[15] A.A. Cooper, T.H. Stevens, Trends Biochem. Sci. 20 (1995) 351–356.
[16] Y. Shao, M.Q. Xu, H. Paulus, Biochemistry 35 (1996) 3810–3815.
[17] M.Q. Xu, D.G. Comb, H. Paulus, C.J. Noren, Y. Shao, F.B. Perler, EMBO J. 13 (1994) 5517–5522.
[18] Y. Shao, M.Q. Xu, H. Paulus, Biochemistry 34 (1995) 10844–10850.
[19] S.F. Altschul, W. Gish, W. Miller, E.W. Myers, D.J. Lipman, J. Mol. Biol. 215 (1990) 403–410.
[20] J.D. Thompson, D.G. Higgins, T.J. Gibson, Nucl. Acids Res. 22 (1994) 4673–4680.
[21] T. Kaneko, S. Sato, H. Kotani, A. Tanaka, E. Asamizu, Y. Nakamura, N. Miyajima, M. Hirosawa, M. Sugiura, S. Sasamoto, T. Kimura, T. Hosouchi, A. Matsuno, A. Muraki, N. Nakazaki, K. Naruo, S. Okumura, S. Shimpo, C. Takeuchi, T. Wada, A. Watanabe, M. Yamada, M. Yasuda, S. Tabata, DNA Res. 3 (1996) 109–136.
[22] A.A. Cooper, Y.J. Chen, M.A. Lindorfer, T.H. Stevens, EMBO J. 12 (1993) 2575–2583.
[23] G. Cullmann, K. Fien, R. Kobayashi, B. Stillman, Mol. Cell. Biol. 15 (1995) 4661–4671.
[24] F.S. Gimble, J. Thorner, Nature 357 (1992) 301–306.
[25] D.A. Shub, H. Goodrich-Blair, Cell 71 (1992) 183–186.
[26] N.F. Neff, Curr. Opin. Cell Biol. 5 (1993) 971–976.
[27] Belfort, M. and Perlman, P.S. (1995) J. Biol. Chem. 270, 30237-20240.
[28] S. Kim, H.G. Dallmann, C.S. McHenry, K.J. Marians, Cell 84 (1996) 643–650.
[29] S. Kim, K.J. Marians, Nucl. Acids Res. 23 (1995) 1374–1379.
[30] M. Reith, J. Munholland, Plant Mol. Biol. Rep. 13 (1995) 333–335.