

## RESEARCH NOTE

## Characterization of an Abundant *Schistosoma mansoni* Transcript with No Homologs in the Databases

Wendell SF Meira, Glória R Franco, Élide ML Rabelo\*, Sérgio DJ Pena/+

Departamento de Bioquímica e Imunologia

\*Departamento de Parasitologia, ICB-UFMG, Av. Antônio Carlos 6627, 31270-010 Belo Horizonte, MG, Brasil

Key words: *Schistosoma mansoni* - expressed sequence tags - database searches - reverse transcribed-polymerase chain reaction

The expressed sequence tag (EST) approach that we have used in the *Schistosoma mansoni* Genome Project is a powerful technique for the discovery of new genes of the parasite (GR Franco et al. 1995 *Gene* 152: 141-147, E Dias Neto et al. 1997 *Gene* 186: 135-142). In a recent comparative study of gene expression in distinct developmental stages of the parasite life cycle using the EST strategy, we identified 466 different genes. From this total, 427 were novel and 333 of them could not be identified based on homologies with database sequences (GR Franco et al. 1997 *DNA Research* 4: 231-240). The high frequency of some of these "unknown" genes in different cDNA libraries suggests that they might have important roles in the biology of *S. mansoni* and thus may constitute possible targets for drug design or vaccine production. One of these genes, highly abundant in one of four adult worm libraries that we are studying in our laboratory, was selected for further characterization.

After clustering analysis of ESTs from different *S. mansoni* cDNA libraries using the program ICATOOLS (Franco et al. 1997 *loc. cit.*), we identified a cluster composed of 16 ESTs from both cDNA ends that corresponded to an unknown gene highly frequent in an adult worm cDNA library. We have called this gene *AUT1* for abundant unknown transcript 1. A single strand consensus of approximately 1.6 kb long was derived from the alignment of the 16 EST sequences using the program DNAsis. In order to obtain the cDNA full-length sequence from both strands, the cDNA clone containing the largest insert was digested with the restriction enzymes *Hind*III and *Sph*I (Fig. 1), the fragments produced were further cloned into the pUC18 vector (Pharmacia) and completely sequenced on both strands using the Thermo-Sequence Fluorescent Labeled Primer Cycle Sequencing kit (Amersham Life Science) and the A.L.F. Automated DNA Sequencer (Pharmacia). Another strategy used to obtain the cDNA sequence from both strands was the amplification by polymerase chain reaction (PCR) of different regions of the cDNA using specific primers designed for the gene. The fragments were cloned into the pUC18 vector (*Sma*I cloning site) using the SureClone Ligation kit (Pharmacia) and sequenced as before. The sequences generated from both processes were aligned using the DNAsis program and the cDNA full-length sequence from both directions was obtained, totaling 1520 bp (Fig. 2). The cDNA was translated into the six possible frames and the length of the longest open reading frame (ORF) was 1005 bp long, potentially encoding a protein of 335 amino acids (Fig. 2).

Analysis of the primary structure of the putative protein coded by *AUT1* gene shows one potential site for N-linked glycosylation, as well as eight sites for phosphorylation by protein kinase C, five for phosphorylation by casein kinase II and eight sites for cAMP-dependent protein kinase, suggesting this protein might be phosphorylated in the organism. The protein does not contain any signal sequence responsible for translocation across the endoplasmic reticulum (ER) membrane, as seen in secretory and plasma membrane spanning proteins, or stretches of hydrophobic residues for plasma membrane insertion. Several searches were performed on distinct databases of protein sequences, typical protein domains and families of proteins as an attempt to identify the predicted protein or a specific domain on it. All the searches provided neither homology with any sequences in databases nor identifiable structural domains.

Primers designed for PCR amplification of a fragment containing the complete coding region

This research was funded by grants from UNDP/WORLD BANK/WHO Special Programme for Research in Tropical Diseases, Fapemig and CNPq/PADCT.

+Corresponding author. Fax: +55-31-441.5963. E-mail: spena@dcc.ufmg.br

Received 4 May 1998

Accepted 31 August 1998

of the cDNA were used for amplification of the *AUT1* gene from the parasite genome. The 1.75 kb fragment obtained was cloned into pUC18 vector. Segments of this insert were amplified using the same primers described before for the amplification of the cDNA. These segments were also subcloned into pUC18 and sequenced in both directions. The generated sequences were aligned using the program DNAsis, given a 1754-bp genomic sequence. The gene contains three introns. The first intron (129 bp) interrupts the coding region of the gene and the other two introns (42 bp and 63 bp) are located towards the 3' untranslated region (Fig. 1). The canonical donor/acceptor splice sites are conserved at all exon/intron junctions.

The expression of the *AUT1* gene in different stages of the parasite life cycle was verified by reverse transcribed PCR (RT-PCR). A 498-bp fragment containing the total 3' untranslated region and the final part of the coding region was amplified

from mRNA preparation from eggs, miracidia, 3h-schistosomula and adult worms. The gene was seen to be expressed in all stages analyzed (data not shown).

In summary, we have characterized an unknown gene of *S. mansoni*, (*AUT1*) that is expressed in different stages of its life cycle. Although the putative protein coded by this gene may potentially be phosphorylated, it does not possess any characteristic structural domain, membrane spanning regions, signals for localization in the nucleus or organelles or signal for translocation across the ER membrane. Homology searches conducted in distinct databases did not show any significant similarity to existing genes, proteins or ESTs. This investigation is being continued in our laboratory. The protein is already expressed in bacteria and polyclonal antibodies are going to be produced for immunolocalization assays, so that we can have some clue of the cellular function of this protein.

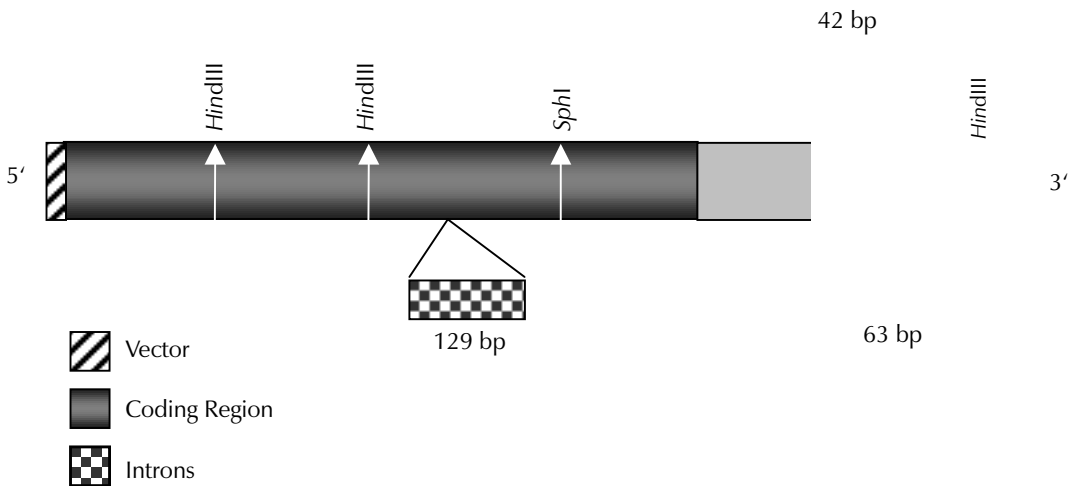


Fig 1: schematic representation of the *AUT1* gene. Positions and sizes of the three introns, the coding region and the restriction sites used for digestion and subcloning are indicated.

1	GGT	GAT	CAT	GTT	GAC	TGC	TTG	TCC	ACA	GTG	TGT	TCT	GTA	GTC	ATG	GCT	GGT	GTT	54
1															<u>M</u>	<u>A</u>	<u>G</u>	<u>V</u>	4
55	GAC	CAC	GCC	TAT	CAG	TAC	TCC	GTG	GAA	GTG	AAG	AGC	AGG	TTC	TCA	CTG	TTC	CTG	108
5	<u>D</u>	<u>H</u>	<u>A</u>	<u>Y</u>	<u>Q</u>	<u>Y</u>	<u>S</u>	<u>V</u>	<u>E</u>	<u>V</u>	<u>K</u>	<u>S</u>	<u>R</u>	<u>F</u>	<u>S</u>	<u>L</u>	<u>F</u>	<u>L</u>	22
109	GAC	GAC	ACT	CTG	AAT	AGT	GAA	GAC	CCA	GAT	ATC	CTG	CTA	TCC	<u>AAG</u>	TTG	CAG	AGT	162
23	<u>D</u>	<u>D</u>	<u>T</u>	<u>L</u>	<u>N</u>	<u>S</u>	<u>E</u>	<u>D</u>	<u>P</u>	<u>D</u>	<u>I</u>	<u>L</u>	<u>L</u>	<u>S</u>	<u>K</u>	<u>L</u>	<u>Q</u>	<u>S</u>	40
163	AAA	CGA	GGT	GAA	AAA	ACT	AAA	AAG	GAT	AAG	CCT	CAT	TTG	CAA	CAG	CAA	CAT	GTG	216
41	<u>K</u>	<u>R</u>	<u>G</u>	<u>E</u>	<u>K</u>	<u>T</u>	<u>K</u>	<u>K</u>	<u>D</u>	<u>K</u>	<u>P</u>	<u>H</u>	<u>L</u>	<u>Q</u>	<u>Q</u>	<u>Q</u>	<u>H</u>	<u>V</u>	58
217	GCA	CCT	ACA	AAG	GCC	GAC	ACG	ATA	ACT	AAG	AGT	GAG	GTG	AAG	CTT	GAT	ACT	GCA	270
59	<u>A</u>	<u>P</u>	<u>T</u>	<u>K</u>	<u>A</u>	<u>D</u>	<u>T</u>	<u>I</u>	<u>T</u>	<u>K</u>	<u>S</u>	<u>E</u>	<u>V</u>	<u>K</u>	<u>L</u>	<u>D</u>	<u>T</u>	<u>A</u>	76
271	ACT	CCG	AAA	GCT	GGT	TCT	CGG	GTG	TCG	AAA	ACC	CCC	AAC	TCC	ACA	GAG	CCG	CCA	324
77	<u>T</u>	<u>P</u>	<u>K</u>	<u>A</u>	<u>G</u>	<u>S</u>	<u>R</u>	<u>V</u>	<u>S</u>	<u>K</u>	<u>T</u>	<u>P</u>	<b><u>N</u></b>	<b><u>S</u></b>	<b><u>T</u></b>	<u>E</u>	<u>P</u>	<u>P</u>	94
325	CCA	GTG	CCA	CCT	GAA	GAC	GTG	CAA	ATT	ACA	TCC	GCT	AAA	GGG	ACA	GAT	GAG	CCA	378
95	<u>P</u>	<u>V</u>	<u>P</u>	<u>P</u>	<u>E</u>	<u>D</u>	<u>V</u>	<u>Q</u>	<u>I</u>	<u>T</u>	<u>S</u>	<u>A</u>	<u>K</u>	<u>G</u>	<u>T</u>	<u>D</u>	<u>E</u>	<u>P</u>	112
379	ATT	AGT	ACC	TTT	NTN	CGT	GGT	CGT	GGC	TCA	GGT	AGG	GGA	ACA	CCT	CGA	GGA	ATG	432
113	<u>I</u>	<u>S</u>	<u>T</u>	<u>F</u>	<u>X</u>	<u>R</u>	<u>G</u>	<u>R</u>	<u>G</u>	<u>S</u>	<u>G</u>	<u>R</u>	<u>G</u>	<u>T</u>	<u>P</u>	<u>R</u>	<u>G</u>	<u>M</u>	130
433	CGT	GTA	GGT	AGA	GGG	CAG	GGG	CCA	CGA	ATA	GCT	CCG	ACA	GAA	GCA	CCT	CAA	GAT	486
131	<u>R</u>	<u>V</u>	<u>G</u>	<u>R</u>	<u>G</u>	<u>Q</u>	<u>G</u>	<u>P</u>	<u>R</u>	<u>I</u>	<u>A</u>	<u>P</u>	<u>T</u>	<u>E</u>	<u>A</u>	<u>P</u>	<u>Q</u>	<u>D</u>	148
487	TCA	GTG	AGT	GAT	TTG	AAT	GCT	CCC	AGA	GGG	TCA	AGC	TTT	GAG	CCT	AGA	GGA	CGT	540
149	<u>S</u>	<u>V</u>	<u>S</u>	<u>D</u>	<u>L</u>	<u>N</u>	<u>A</u>	<u>P</u>	<u>R</u>	<u>G</u>	<u>S</u>	<u>S</u>	<u>F</u>	<u>E</u>	<u>P</u>	<u>R</u>	<u>G</u>	<u>R</u>	166
541	GGG	AGA	GGC	AGA	GGA	AGA	GGA	ATG	TTT	GGT	CGA	GGT	CGT	GGG	ATG	CCA	TTC	AAT	594
167	<u>G</u>	<u>R</u>	<u>G</u>	<u>R</u>	<u>G</u>	<u>R</u>	<u>G</u>	<u>M</u>	<u>F</u>	<u>G</u>	<u>R</u>	<u>G</u>	<u>R</u>	<u>G</u>	<u>M</u>	<u>P</u>	<u>F</u>	<u>N</u>	185
595	TCA	AAT	CGA	GAC	TTC	GAA	AAT	CAG	GAT	GGC	CCA	GAT	CGT	CAG	GGT	CCA	AGA	CAA	648
185	<u>S</u>	<u>N</u>	<u>R</u>	<u>D</u>	<u>F</u>	<u>E</u>	<u>N</u>	<u>Q</u>	<u>D</u>	<u>G</u>	<u>P</u>	<u>D</u>	<u>R</u>	<u>Q</u>	<u>G</u>	<u>P</u>	<u>R</u>	<u>Q</u>	202
649	TAT	GGT	CGA	AGA	GAT	GGG	AAT	TGG	AAT	TCT	CAG	GAT	GTT	GAT	GGT	CTA	ATA	ATG	702
203	<u>Y</u>	<u>G</u>	<u>R</u>	<u>R</u>	<u>D</u>	<u>G</u>	<u>N</u>	<u>W</u>	<u>N</u>	<u>S</u>	<u>Q</u>	<u>D</u>	<u>V</u>	<u>D</u>	<u>G</u>	<u>L</u>	<u>I</u>	<u>M</u>	220
703	CCT	GAG	AGT	GGT	GAC	TCA	GAG	CAA	GTT	GTG	CGT	TTC	GCA	GAC	GAC	CGT	AAC	GAA	756
221	<u>P</u>	<u>E</u>	<u>S</u>	<u>G</u>	<u>D</u>	<u>S</u>	<u>E</u>	<u>Q</u>	<u>V</u>	<u>V</u>	<u>R</u>	<u>F</u>	<u>A</u>	<u>D</u>	<u>D</u>	<u>R</u>	<u>N</u>	<u>E</u>	238
757	GTC	GAG	GAC	CAG	CCA	GAG	CAT	GCC	ACA	GCG	GAG	AAT	GAG	GAG	GGT	GTG	GTT	GTC	810
239	<u>V</u>	<u>E</u>	<u>D</u>	<u>Q</u>	<u>P</u>	<u>Q</u>	<u>H</u>	<u>A</u>	<u>T</u>	<u>A</u>	<u>E</u>	<u>N</u>	<u>E</u>	<u>E</u>	<u>G</u>	<u>V</u>	<u>V</u>	<u>V</u>	256
811	GGC	ACA	GAA	ACC	CCA	GTT	GAG	GAA	GAA	CCG	AAA	AGC	TAC	ACA	CTA	GAA	GGA	TAC	864
257	<u>G</u>	<u>T</u>	<u>E</u>	<u>T</u>	<u>P</u>	<u>V</u>	<u>E</u>	<u>E</u>	<u>E</u>	<u>P</u>	<u>K</u>	<u>S</u>	<u>Y</u>	<u>T</u>	<u>L</u>	<u>E</u>	<u>G</u>	<u>Y</u>	274
865	AAA	GCT	ATG	CGT	CAG	TCT	TCC	AAA	CCA	GCC	GTT	TTA	TTA	AAC	AAC	AAA	GGA	CTA	918
275	<u>K</u>	<u>A</u>	<u>M</u>	<u>R</u>	<u>Q</u>	<u>S</u>	<u>S</u>	<u>K</u>	<u>P</u>	<u>A</u>	<u>V</u>	<u>L</u>	<u>L</u>	<u>N</u>	<u>N</u>	<u>K</u>	<u>G</u>	<u>L</u>	292
919	CGC	AAA	GCT	AAT	GAC	GGC	AAA	GAT	GTG	TTT	GCA	AAC	ATG	GTG	GCT	CAC	AGA	AAA	972
293	<u>R</u>	<u>K</u>	<u>A</u>	<u>N</u>	<u>D</u>	<u>G</u>	<u>K</u>	<u>D</u>	<u>V</u>	<u>F</u>	<u>A</u>	<u>N</u>	<u>M</u>	<u>V</u>	<u>A</u>	<u>H</u>	<u>R</u>	<u>K</u>	310
973	TTA	CAA	GAG	GTC	TCC	GAA	GAT	GTC	TAC	GAG	GTG	GAG	GAG	AGG	AAA	ACG	AGC	CTG	1026
311	<u>L</u>	<u>Q</u>	<u>E</u>	<u>V</u>	<u>S</u>	<u>E</u>	<u>D</u>	<u>V</u>	<u>Y</u>	<u>E</u>	<u>V</u>	<u>E</u>	<u>E</u>	<u>R</u>	<u>K</u>	<u>T</u>	<u>S</u>	<u>L</u>	328
1027	AGA	GCA	TCA	GTC	GAT	CGA	TAT	TGA	TTT	CCC	TTT	CGC	TGA	CGA	CTT	TGG	AAA	CCG	1080
329	<u>R</u>	<u>A</u>	<u>S</u>	<u>V</u>	<u>D</u>	<u>R</u>	<u>Y</u>	<b>***</b>	<b>...</b>	<b>...</b>	<b>...</b>	<b>...</b>	<b>***</b>	<b>...</b>	<b>...</b>	<b>...</b>	<b>...</b>	<b>...</b>	335
1081	GTG	CAG	GGC	CAA	GAA	GTG	AAC	GTG	GGG	GCG	AGG	AAG	AGG	TCA	GAT	AAG	GAG	TGA	1134
1135	CGG	AAG	AGG	GCC	GTG	GTT	TCG	GCC	GTG	GTT	TGC	CTA	TTC	CAA	CGA	GTC	ACA	CTG	1188
1189	CGA	CCT	CCA	GCA	ATT	AAC	AAC	GAC	GCC	AGG	AGT	TCC	CTT	CCT	TGA	AGT	GAA	ATC	1242
1243	TCA	AAG	GAA	AAT	GAA	AAC	AGA	TTG	TTG	GGT	TGA	TTA	TAT	TCT	TCC	TTG	TCA	ACT	1296
1297	TTT	ACC	TTT	CTC	TAA	GGC	TTT	TAG	TCC	AAG	CCT	CTA	CTT	ACT	CAT	ACT	CAT	CTT	1350
1351	TCC	CTA	TTA	AAC	TTT	AGC	TTT	GTA	ATC	TAG	TTT	CTC	CTA	AAA	GCA	GTT	CAA	GGA	1404
1405	TCG	GAT	TTG	GCT	GTG	CTG	CTG	GAG	ATC	TAA	TTG	CAA	AGC	TTT	GTA	CAT	GAG	TGT	1458
1459	TTT	CAG	AAT	GTC	AGT	TTT	GCG	GTT	GAA	GTG	ACC	AAT	CGC	CCC	AAA	TAA	AAA	TTA	1512
1513	ATA	AGA	TG	3'															1520

Fig 2: nucleotide sequence of the cDNA and the putative protein encoded by the *AUT1* gene. The potential sites for protein phosphorylation by PKC - {ST}X{RK}, CKII - {ST}XX{DE} and cAMP-dependent kinase - {RK}XX{ST} are underlined in the sequence. The potential site for N-glycosylation (NX{ST}) is in bold italics.

