

# Cloning and sequencing of a leech homolog to the *Drosophila* engrailed gene

Cathy J. Wedeen, David J. Price\* and David A. Weisblat

Department of Molecular and Cell Biology, University of California, Berkeley, CA 94720, USA

Received 26 November 1990

We have cloned and sequenced a homolog (*ht-en*) to the *Drosophila* engrailed (*en*) gene from the glossiphoniid leech, *Helobdella triseriata*. Amino acid comparisons of the *ht-en* homeobox and C-terminal residues with the corresponding residues encoded by *en*-class genes of other species reveal 75-79% sequence identity. In addition, the *ht-en* sequence appears to have a serine-rich region 16 residues C-terminal from the homeobox, which by analogy to *Drosophila* may be a target site for phosphorylation. The leech gene encodes some amino acid substitutions for residues that are highly conserved in other species. These are found within the second and third of the three putative helices of the homeobox, and in both of the intervening turn regions.

Engrailed: Homeobox gene; *Helobdella triseriata*

## 1. INTRODUCTION

The homeobox gene, engrailed (*en*), encodes a DNA-binding protein that is necessary to establish the 'identity' of the posterior compartment within each segment in *Drosophila* [1-3]. The *en* gene encodes a serine-rich protein that has been shown to be the target of serine phosphorylation [4]; it has been proposed that other segmentation genes (e.g. fused) may regulate *en* function by phosphorylation [5]. A closely related gene in *Drosophila* is *invected* (*inv*), for which no function has yet been determined. Both *en* and *inv* are transcribed concurrently in the same tissues during embryogenesis [6]. In addition, *en* is expressed later in development in certain neurons of the central and peripheral nervous systems [7-10].

*En*-class genes of divergent species are defined as a subfamily of homeobox-containing genes having an especially distinct and highly conserved homeobox region. This high degree of conservation has led to the identification and cloning of homologs from divergent species. In the fruit fly, honeybee, mouse, chicken, zebrafish, and human, two copies of *en*-class genes have been identified; in other species (grasshopper and sea urchin) only one *en*-class gene has been found [10-17]. Thus, it may be that a single *en* gene was present in a common ancestor to the arthropods, echinoderms and chordates and that this gene was

duplicated independently in two, and maybe more, separate lines (i.e. the chordates and the insects).

We have previously reported an *en*-class gene in the leech, *Helobdella triseriata* [18]. We have now cloned and sequenced the homeobox and 3' nucleotides of this gene (*ht-en*) and we compare this sequence with those of other *en*-class genes.

## 2. MATERIALS AND METHODS

### 2.1. Library screening

The phage,  $\lambda$ Ht-en1, was one of 10 recombinants obtained by screening  $6.8 \times 10^4$  plaque forming units from a *Helobdella triseriata* genomic library [19] using the low stringency hybridization conditions described by McGinnis et al. [20]. The probe used was a 250 bp *Pvu*II-*Sau*I *en* cDNA fragment containing the homeobox and upstream region from the clone *en*-HB1 [3].

### 2.2. DNA sequencing

Both strands of the 500 bp *Pvu*II fragment (Fig. 1) were sequenced. Most of the sequence reported in Fig. 2 (i.e. the 3' 147 bp of the homeobox and the downstream region preceding the first termination codon) is a subset of these data. Homeobox sequence 5' to the *Pvu*II site was obtained from a subclone of the 3.5 kb *Hpa*I fragment, using oligonucleotide primers designed to anneal to already sequenced portions of the clone. All sequencing was done using the dideoxy chain termination method.

## 3. RESULTS AND DISCUSSION

### 3.1. The *ht-en* sequence is highly conserved

A recombinant clone homologous to *Drosophila en* was obtained by low stringency hybridization to a *Helobdella triseriata* library (Fig. 1, and section 2). The nucleotide and deduced amino acid sequence of the *ht-en* homeobox and C-terminal flanking region are given in Fig. 2. Given the probe used to clone  $\lambda$ Ht-en1 contained the *Drosophila en* homeobox and 5' sequences,

Correspondence address: C.J. Wedeen, Department of Molecular and Cell Biology, University of California, Berkeley, CA 94720, USA

\*Present address: Department of Physiology, University Medical School, Teviot Place, Edinburgh, EH8 9AG, UK

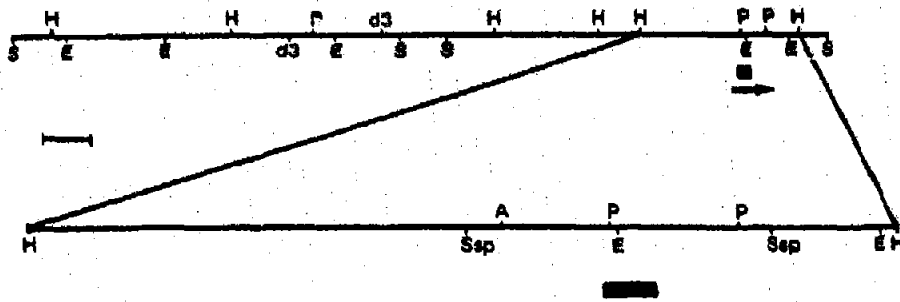


Fig. 1. Restriction map of genomic clone ht-en. The upper line shows the map of a 17 kb fragment. The *Sall* sites at the ends of the clone are from the polylinker of EMBL3 [21]. A blow-up of the 3.5 kb *HpaI* fragment containing the homeobox is shown on the lower line. The position of the homeobox is shown by the filled boxes below each line. The arrow below the upper line designates the putative direction of transcription. The scale bar is equivalent to 1 kb for the upper line and 200 bp for the lower line. Key: A, *ApoI*; d3, *HindIII*; E, *EcoRI*; H, *HpaI*; P, *PvuII*; S, *Sall*; Ssp, *SspI*.

it was expected that the homeodomain portion of the cloned leech gene should be homologous to *en*. We do indeed observe this expected homology, but in addition, there is extensive homology extending 19 residues C-terminal to the homeobox, in a region not represented by the probe (Fig. 3). By these criteria we designate ht-en as an *en* homolog. The inferred amino acid sequence of the entire conserved region of ht-en was compared to the corresponding region of the other *en*-class genes from the species listed in Fig. 3. The ht-en amino acid sequence is 75-79% identical to the other *en* homologs. Given the short sequence length, the high degree and close range of sequence identity, we are unable to make any significant correlations between the degree of sequence identity and the time since evolutionary divergence of the leech from any of these species.

3.2. *ht-en* encodes amino acid substitutions in the homeodomain

X-Ray diffraction has been used to determine the structure of an *en* homeodomain/DNA complex [23]. The proposed structure of the *en* homeodomain is similar to that proposed for the *Antennapedia* homeodomain on the basis of nuclear magnetic resonance [24]. The *en* homeodomain contains 3  $\alpha$ -helices and an N-terminal arm. Helices 1 and 2 pack against each other in an antiparallel arrangement and make few contacts with the DNA; helix 3 lies perpendicular to helices 1 and 2 and, as the 'recognition helix', makes extensive contacts with the major groove of the DNA. The residues composing each of the helices are designated in Fig. 3. In the ht-en homeodomain several amino acid changes are observed. Some of these amino acid differences have been reported earlier in a discussion of the epitope for a monoclonal antibody, mab4D9, directed against a portion of the injected homeodomain [10]. Here we describe the substitutions in the ht-en protein with respect to the proposed homeodomain structure. One change occurs at residue 58 within the 'recognition helix', number 3. This residue is isoleucine in every *en*-class protein except *Drosophila inv*, where it is leucine, and in the ht-en homeodomain, where it is a methionine. The other changes occur in helix 2 and in the turn regions between helices 1 and 2, and between helices 2 and 3. Substitutions within helix 2 occur at residues 34 and 35. One or both of these is always glutamine except in the sea urchin, where they are arginine and serine, and in leech, where they are threonine and cystine. In the turn regions, residue 26 is always arginine except in the sea urchin, where it is asparagine, and in leech, where it is lysine; residue 41 is often glycine in *en*-class genes but in sea urchin and in mouse this residue is replaced by the more sterically restricting threonine and serine, respectively, and in leech an asparagine, a nonconservative substitution, is found in this position. None of these amino acid substitutions occurs in a position that has

1	CAG GAC GAA AAG AGA CCT GGA ACA GCA TTC ACG GGC GAT CAG CTG GCG	48
	GLN ASP GLU LYS ARG PRO ARG THR ALA PHE THR GLY ASP GLN LEU ALA	
49		10
	AGG TTG AAG CGT GAA TTC AGC GAG AAC AAA TAC CTG ACG GAG CAG AGG	56
	ARG LEU LYS ARG GLU PHE SER GLU ASN LYS TYR LEU THR GLU GLN ARG	
97		20
	AGA ACA TGT CTG GCG AAG GAA CTG AAC TTG AAC GAG AGC CAG ATC AAA	144
	ARG THR CYS LEU ALA LYS GLU LEU ASN LEU ASN GLU SER GLN ILE LYS	
145		40
	ATC TGG TTC CAG AAC AAG AGG GCC AAG ATG AAC AAG GCG AGT GGC GTG	192
	ILE TRP PHE GLN ASN LYS ARG ALA LYS MET LYS LYS ALA SER GLY VAL	
193		50
	AAG AAT CAG TTG GCT CTG CAA CTC ATG GCA CAG GGC CTC TAC AAC CAC	240
	LYS ASN GLN LEU ALA LEU GLN LEU MET ALA GLN GLY LEU TYR ASN HIS	
241		70
	TCA TCA TCA TCA TCT TCT TCT TCC TCC TCC TCC TCT TCG ATC TTC CTC	280
	SER SER SER SER SER SER SER SER SER SER SER SER SER ILE PHE LEU	
290		80
	CTC GCA TAA	280
	LEU ALA	280
98		90

Fig. 2. Nucleotide and deduced amino acid sequence of the ht-en homeobox and 3' flanking region. The 294 nucleotide sequence of the putative open reading frame containing the homeobox and 3' sequences and ending in a stop codon, TAA, is shown on the upper line; the homeobox (nucleotides 3-186) is underlined. The first and last nucleotides of each line are numbered above the line. The corresponding amino acid sequence of the putative open reading frame is given on the lower line; every 10 amino acids are numbered below the line.

