# Summary

**Thesis summary**

I have taken two complementary approaches to isolating cell-type specific *cis*-regulatory regions upstream of three genes, *egl-17*, *zmp-1* and *cdh-3*.   In the first approach (Chapter 2), I used a sufficiency analysis to test genomic regions of DNA upstream of three genes for their ability to confer cell-specific expression on a naïve promoter, *pes*-10.  In a second, orthogonal, approach (Chapter 3), I compared homologous upstream regions (phylogenetic footprints) to identify regions of similarity responsible for conferring cell type-specific patterns of expression.

The selection of these three genes stemmed from the fact that they are expressed in a restricted number of overlapping cell types at similar times. Genes that are specifically expressed in the same tissue at the same time might have common regulatory programs and might be recognized by common *trans* factors. Therefore, conserved motifs in genes showing common expression profiles are likely to be involved in spatial/ temporal expression. Additionally, with the exception of the early expression of *egl-17* in the presumptive vulE and vulF cells, all vulval and anchor cell expression occurs after terminal differentiation. The isolation of elements that drive post-terminal differentiation expression allows us to determine what makes each of these cell types unique, and to try to make connections between the known signaling pathways involved in these cell's specification and terminal fates decisions.

While it seems that no single approach is going to identify and define all the *cis*-acting regulatory elements responsible for conferring cell type-specific expression, the corroboration of approaches allows for significant progress to be made.

**Sufficiency analysis**

The goals of this study was to define the minimal sequences responsible for conferring specificity off a naïve promoter to several vulval cells and the anchor cell in order to search the genome for similar elements. I have narrowed down a 3.9 kb region to: a 143 bp region of *egl-17* that drives vulC and vulD expression, and a separate 102 bp region that is sufficient to drive the early expression in presumptive vulE and vulF cells. I have narrowed a 3.5 kb region to a 300bp region of *zmp-1* that is sufficient to confer expression in vulE, vulA and the anchor cell. And finally, I have examined a 6.0 kb region to define a 689 bp region of *cdh-3* that is sufficient to drive expression in the anchor cell and vulE, vulF, vulD and vulC; a 155 bp region that is sufficient to drive anchor cell expression; and a separate 563 bp region that is also sufficient to drive expression in these vulval cells. One theme that remains the same in all three analyses is that I failed to identify any repressor elements involved in conferring expression in terminally differentiated cell types. Furthermore, it became clear from this study that there are multiple mechanisms used to ensure fidelity of expression patterns even between genes that are expressed in the same cell. These mechanisms include: the use of discrete separable elements that confer cell-type specific expression (*cdh-3* anchor cell expression and *egl-17* expression in sister cells vulC and vulD); the use of complex patterns of binding sites that combinatorially act to establish the fidelity of expression in a variety of cell types from different lineages (*zmp-1*); and the use of tissue-specific elements responsible for driving expression in an entire tissue rather then in sub-domains of its constituent cells (*cdh-3*).

**Determining the necessity of regions defined by sufficiency analysis**

In one sense, the necessity of these elements was irrelevant to our immediate goal of determining sequences that possess the ability to confer cell type-specific expression of these genes. In our case, the genes themselves are somewhat superfluous compared to the elements, which are sufficient to confer this specificity. What the necessity testing will be invaluable for is putting the results of these analyses back into the context of the native promoters. It will be especially interesting to observe the relative importance of the two non-overlapping regions in upstream sequences of *cdh-3*, both of which, despite qualitative differences, appear sufficient to confer expression in the same cells. Additionally, mutation analysis of the individual elements defined in the sufficiency and phylogenetic footprint studies will allow us to further delimit the boundaries of these regions. If conducted in the context of the native promoter, the significance of these mutations may be weighed in the natural milieu of the gene.

**Phylogenetic footprinting studies of *cis*-regulatory sequences**

Since continuously occurring mutational events accumulate at neutral positions but are eliminated in functional regions, it is argued that conserved motifs in diverse orthologous promoter sequences are more likely to have a functional role (Tagle *et al.*, 1988). In this study, I used two species of *Caenorhabditis*, *C. elegans* and *C. briggsae*, for sequence comparisons. With a two-species comparison, I was able to identify several blocks of homology. In the cases of *zmp-1* and *cdh-3*, these blocks were located throughout the upstream region, and only by using the sufficiency data was I able to hone in on a single

block in each as conferring expression in the anchor cell and/or the vulva cells.

Presumably, these other blocks of similarities throughout the upstream regions confer

expression in other cell types, as these markers are expressed in a variety of tissues. In the

case of *egl-17*, the only elements found, by our sequence comparison were in a region

that was found to drive expression in the vulval cells.  This is not surprising since the

expression of this marker is restricted to very few tissues.

The regions of similarity that did direct vulval and anchor cell specific expression

are still broad enough to obscure the resolution of distinct binding sites; furthermore,

multiple *trans*-acting sites may be needed to confer a specific expression pattern. In order

to get a more defined picture of the regions I have found, it will be helpful to compare co-

regulated or homologous genes from several other species in order to distinguish signal

from the background noise. With the addition of other species, it may be possible to

define this region in greater detail. The present nematode tree gives two additional

siblings, CB5161 and PS1010, that may be very useful for such comparisons (Figure 1)

(Fitch et al., 1995). I am currently trying to isolate the upstream regions of the *egl-17*,

*zmp-1* and *cdh-3* genes from these species for use in a four-way comparison. As one adds

more species to the analysis, the distinction between conserved motif and diverged

background should become clearer. One risk with this type of analysis is that when

including many sequences, particularly distantly related ones, there is an increased

chance that some of them may have lost, or completely altered, some regulatory elements

over the course of evolution (reviewed in Blanchette and Tompa, 2002). This makes the

selection of species imperative to the successful outcome of the analysis. One advantage

of this type of approach over others is that while other approaches will distinguish a

single site as necessary and/or sufficient, this approach may help delimit multiple

elements in the *cis*-acting regions to give a broader view of the *cis*-acting sequences.

**Practical considerations when identifying phylogenetic footprints**

ClustalW (Higgins *et al.*, 1994) alignments do not always work for identifying such

footprints. Regulatory elements tend to be short (8-10 bp) relative to the entire regulatory

region. If the species are more diverged, the noise of the diverged nonfunctional

background will overcome the short conserved signal. The result is that the alignment

will not align the short regulatory elements well; the regulatory elements would go

undetected. This failure was the case for the *zmp-1* and *cdh-3* upstream regulatory

regions. There is enough divergence in the sequence that elements picked up by the

Seqcomp and Family Relations programs were completely obscured in the clustalW

alignment (data not shown). The *egl-17* clustalW alignments were able to identify regions

of similarity (data not shown). However, there are large blocks of similarity in the

upstream sequences of this gene, making this method, while still fruitful, less helpful than

in the case of the other two genes. Additionally, many alignment tools and comparisons

do not allow the identification of reverse complement similarities, which can be

functionally significant in the context of enhancers that may operate in either direction.

**Combining the results of sufficiency testing and phylogenetic footprinting studies**

By combining the results of my sufficiency testing with the results of the phylogenetic

footprinting, it was satisfying to find that both methods were able to hone in on similar

regions as those that were important for conferring tissue specific expression. As can be

seen in figures 2 (*egl-17*), 3 (*zmp-1*), and 4 (*cdh-3*), there are conserved elements that fall

in the regions of sufficiency in each of these three genes. In the case of *egl-17*, the

location of element D, which falls in the middle of the minimal region defined by

sufficiency, is very encouraging; putative binding sites or over-represented sequences in

this region should provide good candidates for cell-specific elements. The location of

element B in *egl-17* is in a region that plays a role in conferring GFP expression in vulE

and vulF.  In *zmp-1*, the locations of all conserved elements appear to fall in regions that

were important for vulE, vulA and anchor cell expression.  Multiple conserved elements

in *cdh-3* are found in the regions defined by sufficiency analysis to be important for

vulval and anchor cell expression.


**Analysis of putative *trans*-acting factors**

The sufficiency analysis and phylogenetic footprinting experiments defined overlapping

regions of importance in conferring cell-type specific expression of several vulva cells

and the uterine anchor cell. However, these regions are still broad enough to obscure the

resolution of distinct binding sites. To identify putative *trans*-acting factors that drive

expression in these cells, I turned to the Transfac database (see Transfac analysis in

Chapters 2 and 3) and our knowledge of genes that are likely to be involved in the

specification of these cells (Table 1).

In *lin-29* animals (Horvitz et al., 1983), a gene involved in the heterochronic

pathway (Arasu et al., 1991; Bettinger et al., 1997), *egl-17* expression in the presumptive

vulE and vulF cells persists, and vulC and vulD expression does not ensue. In the case of

*zmp-1*, there is no vulE expression in the young adult, though this background does not

affect the *cdh-3* expression during the L3 and early L4 stages. Since this mutation causes the reiteration of earlier developmental stages, it is not surprising that early expression persists at the expense of the later expression pattern (M. Wang and T. Inoue, unpublished observations).

In the PAX family member *egl-38* (Chamberlin et al., 1997), there is no *egl-17* expression in the presumptive vulE and vulF cells, and no *zmp-1* expression in vulE. The HOM-C family member *lin-39* also decreases the *egl-17* expression in the presumptive vulE and vulF cells, suggesting that these genes may play a role in regulating expression in vulE (M. Wang, unpublished observations)

In animals mutant in the *lin-1* gene (Beitel et al., 1995), which encodes an ETS family member, there is no *egl-17* expression in the presumptive vulE and vulF cells. However, *zmp-1* expression in vulE is normal in the *lin-1* background. *lin-1* also effects vulC and vulD expression in the *egl-17* background. This altered expression suggests that *lin-1* may play a specific role in *egl-17* regulation (M. Wang, unpublished observations).

In the *lin-26* animals, a predicted zinc finger transcription factor that plays a role in the generation of Pn.p cells (Labouesse et al., 1994), *egl-17* vulC expression is lost and the vulD expression is dramatically reduced. Additionally, *cdh-3* expression is dramatically reduced in vulC, D and E (T. Inuoe, unpublished observations). The *lin-26* gene may play an important role in the specification of these cells.

In animals carrying one allele of the gene encoding a GTX NKx6.2 family member *cog-1* (R. Palmer *et al.,* in press), *sy275*, *egl-17* vulE expression is seen in addition to vulC and vulD expression in the L4 stage. This expression is separate from

the early expression in this cell. This same allele shows no vulE *zmp-1* expression.

Perhap, *cog-1* (*sy275*) plays a role regulating late *egl-17* expression in vulE cells (M.

Wang and T. Inoue, unpublished observations). However, no GTX binding sites were

found using the MatInspector program. A second allele of *cog-1*, *sy607*, does not effect

vulE expression. However, this allele shows no *cdh-3* expression in vulC and vulD cells,

and a dramatic reduction in vulE cells (M. Wang and T. Inoue, unpublished

observations).

In the LIM domain protein, *lin-11* (Freyd et al., 1990), there is no *egl-17*

expression in vulC or vulD cells, but there is no effect on the early expression in the

presumptive vulE and vulF cells. In *lin-11*, there is also no *zmp-1* expression in either

vulA or vulE cells, yet it also alters *cdh-3* expression levels in vulF, vulE, vulC and vulD.

This result is surprising because of the *lin-11* effect on *zmp-1* and *cdh-3* expression in the

primary lineage. Although we know that *lin-11* animals have altered secondary cell

lineage, we have no evidence of it having any effects on the analysis of primary fate (B.

Gupta, unpublished observations). Our analysis using the MatInspector program did

identify binding sites for the putative LIM homolog, ISLI-1, in conserved regions

responsible for driving *egl-17* expression in vulC and vulD. The significance of this

finding is not known. This site came up in all the analyses, and has a very loose

consensus sequence with a core matrix sequence of TAAT similar to that of other

homeodomains.

A loss of function mutation in *lin-17* (Sternberg and Horvitz, 1988), which

encodes a WNT-family receptor, causes variable *cdh-3* expression in vulC and vulD and

ectopic variable expression in vulA and vulB (T. Inoue, unpublished observations). This

result suggests that this gene probably plays an intimate role in mediated secondary cell fate or transcriptional regulation.

An anchor cell element that drives transcription of LIN-3 has been isolated, and involves *trans*-acting factors that bind to a nuclear hormone receptor site and E-box protein-binding sites (B. Hwang and P. Sternberg, unpublished results). Disruption of these elements does not disrupt the expression of *cdh-3* or *zmp-1::gfp* in the anchor cell. A different mechanism and/or factors must be used to establish the anchor cell expression of these late markers. We have few candidate factors that may be involved in the regulation in this cell.

While the focus of this project was to isolate cell-specific response *cis*-regulatory elements rather than identifying *trans*-acting factors, I was also looking forward to the more distant goal of determining the integration of signaling pathways in the downstream targets of these pathways. The integration, in the upstream sequences, of members of the RAS, NOTCH and WNT pathways, whose signaling is intimately bound with the establishment of these fates, would help establish the hierarchy of action of these pathways and their interactions. In the case of the early expression of the *egl-17* gene (expression in the presumptive vulE and vulF cells), it is still a matter of debate regarding the determination status of these cells at the time of this expression. *egl-17* is expressed at a time when crucial signaling events that result in an invariant cell fate pattern are still occurring, which makes this particular gene, and the elements responsible for conferring its early expression, of special interest. There are several approaches to the identification of the *trans*-acting factors involved in conferring the cell type-specific expression patterns. The preceding section has talked about various genetic backgrounds that have

been examined in the context of the full-length reporter constructs. Some of these genetic

backgrounds have a dramatic effect on the ability of these reporters to confer expression.

One approach is to use the minimal sufficiency regions defined in this thesis to look at

the genetic backgrounds that had an effect on expression patterns, to establish that they

are working through these elements, and also to extend this to a greater diversity of

genetic backgrounds. This, however, will not get to the crux of the matter of whether

these factors are directly binding these sequences, or are regulating something in turn that

is directly binding them. It will, however, tell you which genes appear to be involved in

establishing the differential gene expression in these cells.

To categorically establish which *trans*-acting factors are binding these sites

directly will require biochemical testing of the ability of a specific *trans*-acting factor to

bind a particular sequence.

**Genomic analysis**

Once elements responsible for conferring cell-type specific expression have been defined

as concisely as bench-work will allow us (through mutational analysis, or further

phylogenetic analysis), it will be both feasible and exciting to search the genome of *C.*

*elegans* and *C. briggsae* for other genes whose *cis*-regulatory sequences contain these

elements.

When a single promoter sequence is searched, one often finds many putative

elements conserved all over the sequence, making it difficult to choose for further

experimental analysis. On the other hand, when multiple promoter sequences are

searched simultaneously, the conserved motifs are more likely to be functionally

important. To this end, I used the AlignACE program to look for over-represented

sequences in elements of intergenic regions found in our sufficiency analysis; I also

looked for over-represented sequences between elements that conferred the same cell

specificity (Chapter 2, Table 2 and Chapter 3, Table 3). One caveat of this approach is

that its efficacy, while seemingly good in yeast (Hughes *et al.*, 2000), has not been tested

on metazoans. The metazoans have much larger non-coding regions use a more

combinatorial based system of regulation show long distance regulation via chromatin,

and appear to have a vast number of transcription factors not present in yeast.  These

over-represented sequences that fall into regions which, by our other analysis, appear to

be important in conferring cell/ tissue specificity make good candidates to search for in

the genome, and also make good candidates for mutational analysis.  In order to perform

this search with a consensus sequence, we can modify the program ScanACE, which

performs a similar search on the genome of *Saccharomyces cerevisiae*.

**References**

Arasu, P., Wightman, B., and Ruvkun, G. (1991). Temporal regulation of *lin-14* by the
     antagonistic action of two other heterochronic genes, *lin-4* and *lin-28*. *Genes &
     Development* **5,** 1825-1833.

Beitel, G., Tuck, S., Greenwald, I., and Horvitz, H. (1995). The *Caenorhabditis elegans*
     gene *lin-1* encodes an ETS-domain protein and defines a branch of the vulval
     induction pathway. *Genes & Development* **9,** 3149-3162.

Bettinger, J., Euling, S., and Rougvie, A. (1997). The terminal differentiation factor LIN-
     29 is required for proper vulval morphogenesis and egg laying in *Caenorhabditis
     elegans*. *Development* **124,** 4333-4342.

Blanchette, M., and Tompa, M. (2002). Discovery of regulatory elements by a computational method for phylogenetic footprinting. *Genome Research* **12,** 739-748.

Chamberlin, H., Palmer, R., Newman, A., Sternberg, P., Baillie, D., and Thomas, J. (1997). The PAX gene *egl-38* mediates developmental patterning in *Caenorhabditis elegans*. *Development* **124,** 3919-3928.

Fitch, D., Bugaj-Gaweda, B., and Emmons, S. (1995). 18S Ribosomal RNA gene phylogeny for some *Rhabditidae* related to *Caenorhabditis*. *Molecular Biology and Evolution* **12,** 346-358.

Freyd, G., Kim, S., and Horvitz, H. (1990). Novel cysteine-rich motif and homeodomain in the product of the *Caenorhabditis elegans* cell lineage gene *lin-11*. *Nature* **344,** 876-879.

Higgins, D.G., Thompson, J.D., and Gibson, T.J. (1988). Using ClustalW for multiple alignments. *Methods Enzymol.* **266**, 387-402.

Horvitz, H., Sternberg, P., Greenwald, I., Fixsen, W., and Ellis, H. (1983). Mutations that affect neural cell lineages and cell fates during the development of the nematode *C. elegans*. *Cold Spring Harbor Symposia on Quantitative Biology* **48,** 453-463.

Hughes, JD., Estep, PW., Tavazoie, S., and Church, GM. (2000). Computational identification of *cis*-regulatory elements associated with groups of functionally related genes in *Saccharomyces cerevisiae*. *Journal of Molecular Biology* **296**, 1205-1214.

Labouesse, M., Sookhare, S., and Horvitz, H. (1994). The *Caenorhabditis elegans* gene *lin-26* is required to specify the fates of hypodermal cells and encodes a presumptive zinc-finger transcription factor. *Development* **120,** 2359-2368.

Sternberg, P., and Horvitz, H. (1988). *lin-17* mutations of *Caenorhabditis elegans* disrupt certain asymmetric cell divisions. *Developmental Biology* **130,** 67-73.

Tagle, D., Koop, B., Goodman, M., Slightom, J., Hess, D., and Jones, R. (1988). Embryonic epsilon and gamma globin genes of a prosimian primate (*Galago crassicaudatus*). Nucleotide and amino acid sequences, developmental regulation and phylogenetic footprints. *Journal of Molecular Biology* **3,** 439-455.

**Figure 1: Selection of nematode species for comparative genomic analysis**

Closest sibling species to *C. elegans* are listed in this tree diagram, adapted from Fitch *et al.*, 1995. Dates of divergence are hard to predict, but the current prediction of divergence between *C. elegans* and *C. briggsae* is 50-120 million years.
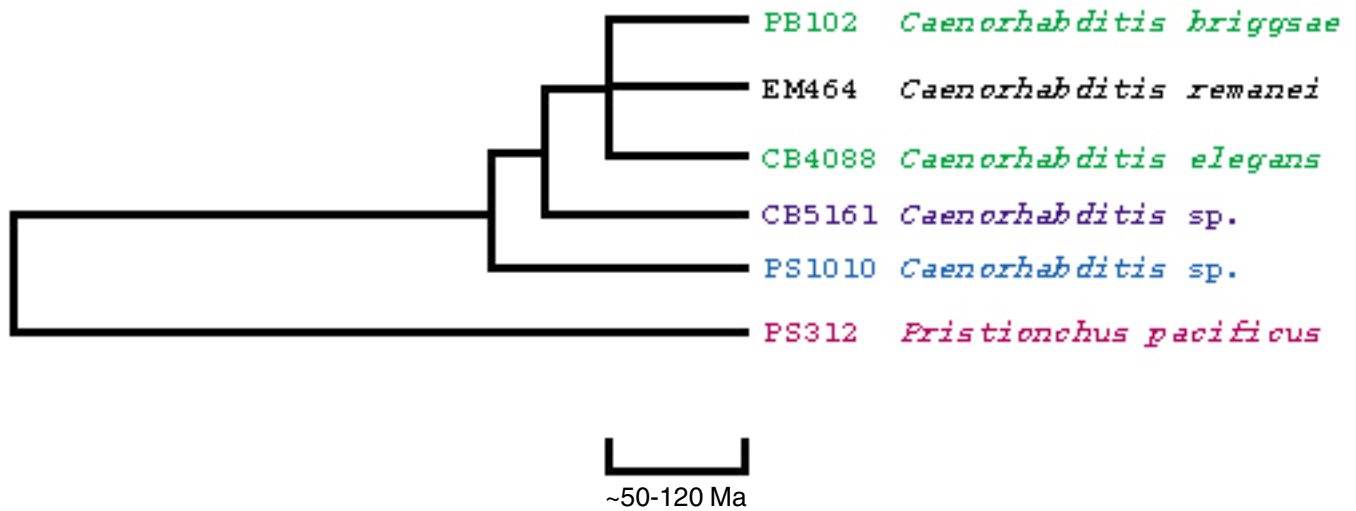
**Figure 1: Selection of nematode species for comparative genomic analysis**



~50-120 Ma

**Figure 2: Combined results of the *egl-17* sufficiency and phylogenetic analyses**

This figures depicts both the *egl-17* sufficiency data as seen in Chapter 2, Figure 3, and

the conserved regions identified in the phylogenetic footprinting studies, which have been

superimposed on this schematic. A, B, C, D represents element A, element B, and so

forth. The boundaries of each element are listed in the top right-hand corner of the figure.

The box in the upper right-hand corner depicts the expression pattern of each of the three

markers used in these studies.

**Figure 2: Upstream regions that direct *egl-17* expression**

*egl-17*  *zmp-1*  *cdh-3*

vulA
vulB1
vulB2
vulC
vulD
vulE
vulF
AC

4270  4298    4466  4485
C    elem D    E
4363    4459

4331  4359    4466  4474    4732  **Expression**

Construct    vulC  vulD

mk80-132 (158 bp)    4316    4474    +    +

mk80-104 (150 bp)    4316    4466    +/-    –

mk125-132 (143 bp)    4331    4474    +    +

mk102-56 (157 bp)    4359    4516    –    +/-

mk102-104 (107 bp)    4359    4466    –    –

B
3208  3241

4708  start *egl-17* txn

1716    3182    3611    4516    4732

Construct    **Early Exp.**

mk15-20 (1974 bp)    1716    3690    +/-

mk82-100 (723 bp)    2888    3611    +/-

mk84-20 (508 bp)    3182    3690    +/-

mk82-85 (315 bp)    2888    3203    -

mk27-20 (188 bp)    3502    3690    –

mk15-148 (3016 bp)    1716    4732    +

mk84-148 (1550 bp)    3182    4732    +

mk103-148 (305 bp)    4427    4732    +/-

mk103-56 (89 bp)    4427    4516    –

mk153-148 (167 bp)    4565    4732    –

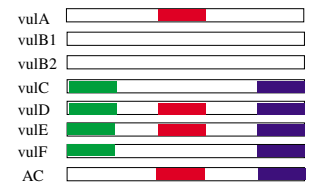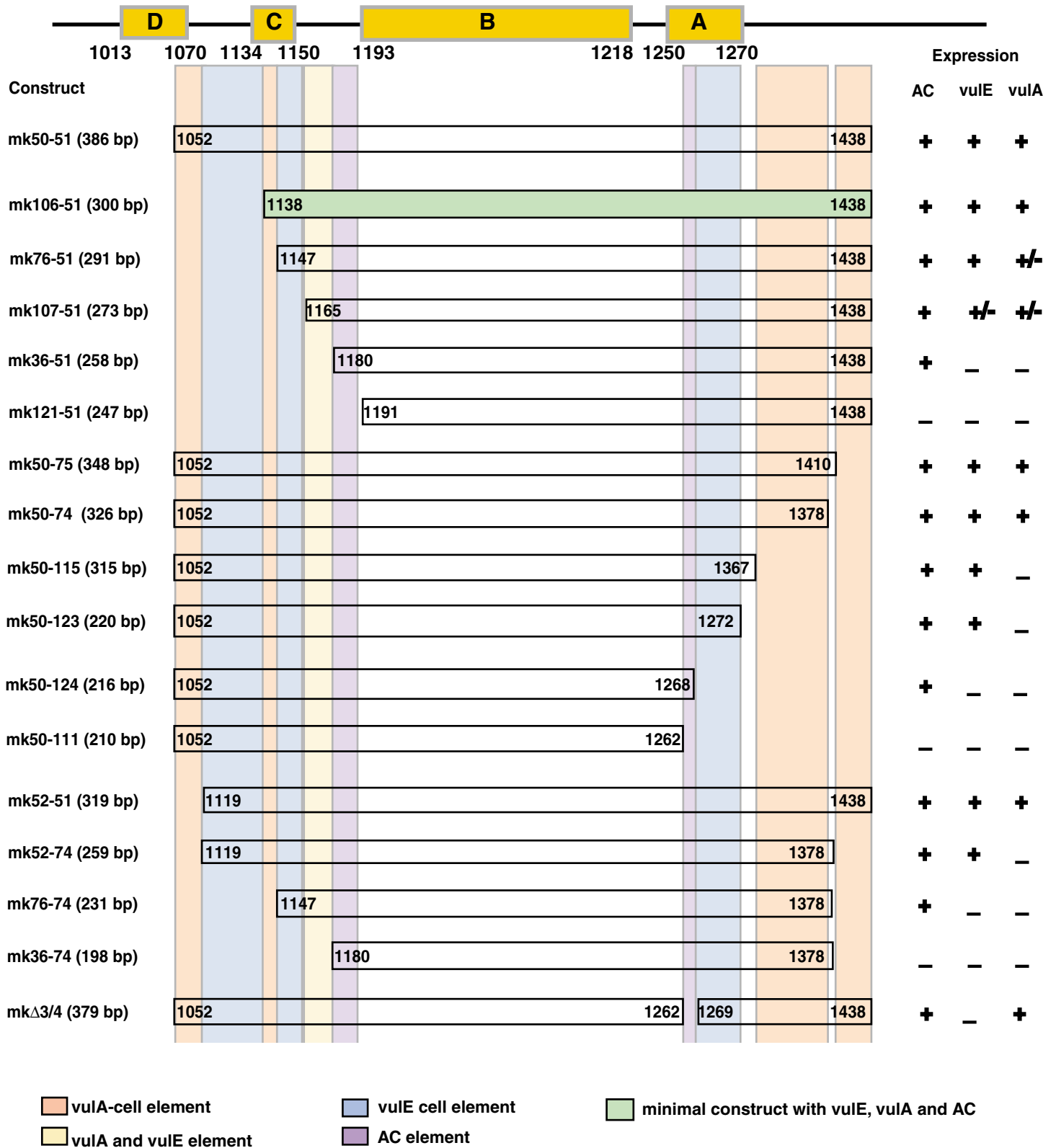mk153-154 (102 bp)    4565    4667    +

Figure 3: Multiple Regions Direct *zmp-1* expression

**Figure 3: Combined results of the *zmp-1* sufficiency and phylogenetic analyses**

This figure depicts both the *zmp-1* sufficiency data as seen in Chapter 2, Figure 5, and the conserved regions identified in the phylogenetic footprinting studies, which have been superimposed on this schematic. A, B, C, D represents element A, element B and so forth. The boundaries of each element are indicated at the bottom of each element. The box in the upper right-hand corner depicts the expression pattern of each of the three markers used in these studies.

**Figure 4: Combined results of the *cdh-3* sufficiency and phylogenetic analyses**

This figure depicts both the *cdh-3* sufficiency data as seen in Chapter 2, Figure 6, and the conserved regions identified in the phylogenetic footprinting studies, which have been superimposed on this schematic. A, B, C, D represents element A, element B and so forth. Elements H, I, J, K are overlapping a consecutive, and so have been represented by a single box labeled "HIJK". The boundaries of each element are listed in the top right-hand corner of the figure. The box in the upper right-hand corner depicts the expression pattern of each of the three markers used in these studies.
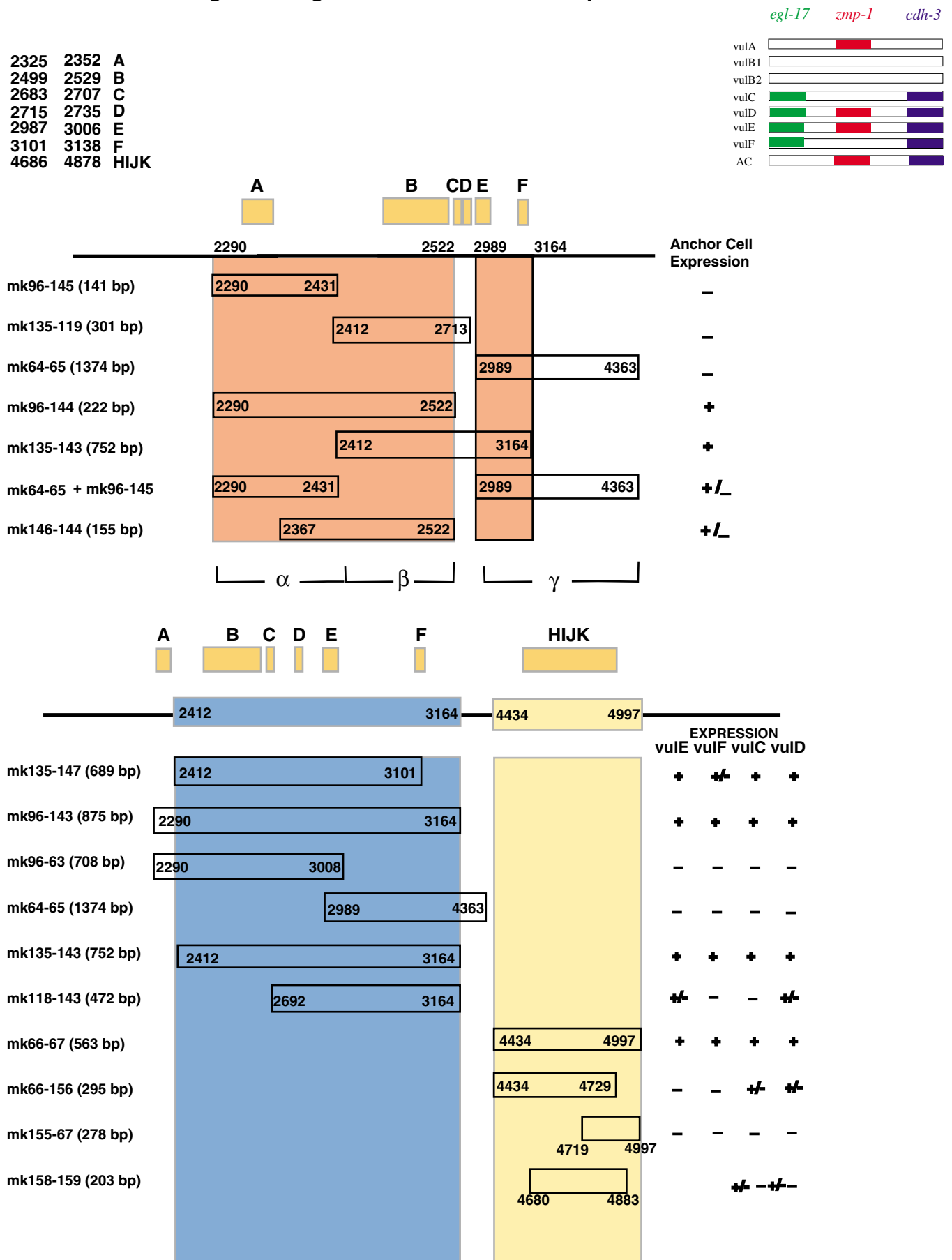
# Figure 4:Regions that direct *cdh-3* Expression



Figure 4: Regions that direct *cdh-3* Expression

**Table 1: Effect of genetic background on marker expression**

For each marker gene listed in the first column, the expression pattern in a variety of different genetic backgrounds (listed in column two) is summarized for cells vulA-F. The expression pattern in the anchor cell was not determined. An "nd" means that the expression pattern was not determined. A "+/-" indicates that expression was variable or weak. (These data summarize expression studies preformed by M. Wang and T. Inoue, unpublished results.)

**Table 1: Effect of genetic background on marker expression**

| marker | Genetic background | vulA | vulB | vulC | vulD | vulE | vulF |
|---|---|---|---|---|---|---|---|
| *egl-17::*GFP | wt | | | + | + | + | + |
| | *lin-29 (sy292 /n333)* | | | - | - | + (persists longer) | + (persists longer) |
| | *lin-26 (ga91)* | | | - | +/- | nd | nd |
| | *cog-1 (sy275)* | | | + | + | ++ (at time 2°) | + |
| | *cog-1 (sy607)* | | | + | + | + | + |
| | *lin-11 (n389)* | | | - | - | + | + |
| | *egl-38 (n578)* | | | nd | nd | - | - |
| | *lin-1 (sy254)* | | | +/- | +/- | - | - |
| | *lin-39 (n709)* | | | nd | nd | +/- | +/- |
| | *lin-17* | | | nd | nd | nd | nd |
| | *sqv-3 (n2842)* | | | + | + | + | + |
| | *evl-2 (ar101)* | | | + | + | + | + |
| | *evl-22 (ar104)* | | | + | + | + | + |
| | | | | | | | |
| *cdh-3::*GFP | wt | | | + | + | + | + |
| | *lin-29* | | | + | + | + | + |
| | *lin-26 (ga91)* | | | +/- | +/- | +/- | + |
| | *cog-1 (sy275)* | | | + | + | + | + |
| | *cog-1 (sy607)* | | | - | - | +/- | + |
| | *lin-11* | | | - | - | - | +/- |
| | *egl-38 (n578)* | | | nd | nd | nd | nd |
| | *lin-1* | | | nd | nd | nd | nd |
| | *lin-39* | | | nd | nd | nd | nd |
| | *lin-17* | +/- | +/- | +/- | +/- | + | + |
| | | | | | | | |
| *zmp-1*::GFP | wt | + | | | | + | |
| | *lin-29 (sy292)* | + | | | | - | |
| | *lin-26 (ga91)* | nd | | | | nd | |
| | *cog-1 (sy275)* | + | | | | - | |
| | *cog-1 (sy607)* | nd | | | | nd | |
| | *lin-11 (n389)* | - | | | | - | |
| | *egl-38 (n578)* | nd | | | | - | |
| | *lin-1 (sy254)* | nd | | | | + | |
| | *lin-39* | nd | | | | nd | |
| | *lin-17* | nd | | | | nd | |
| | *lin-31 (n301)* | + | | | | nd | |