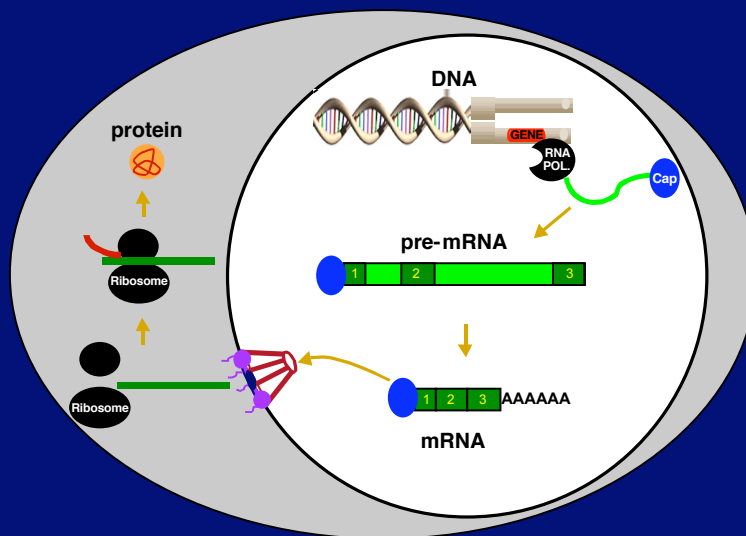# Computational Molecular Biology of Genome Expression and Regulation

**Michael Q. Zhang**
**Cold Spring Harbor Labratory**

- An ESE (enhancer-sequence-element) SNP can alter RNA splicing (Brac1:exon18)
- Classification of 5'UTRs by *CART*
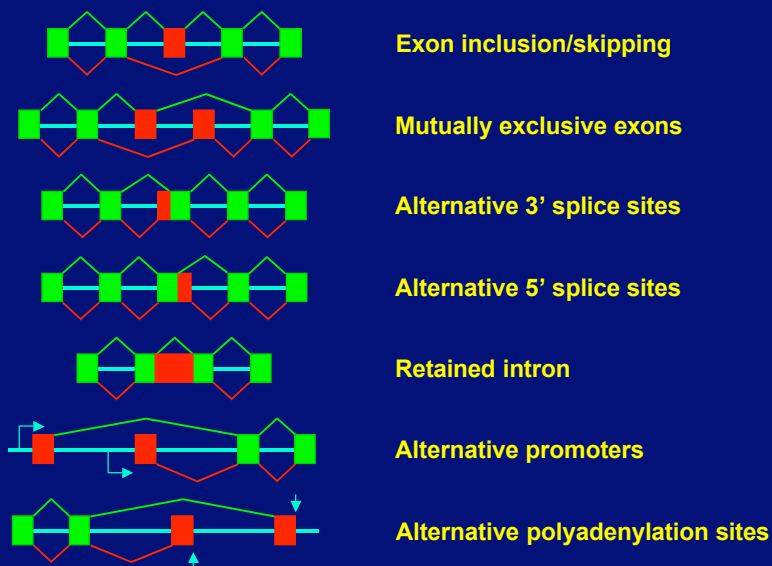- Promoter prediction and analysis
- *CSEdb* and human gene number
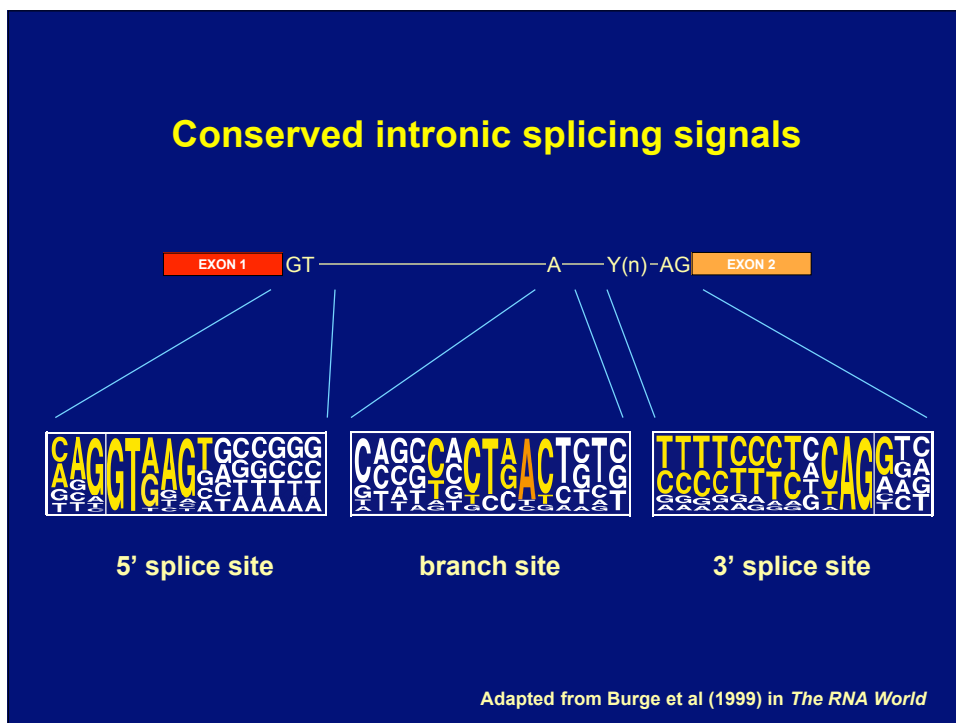
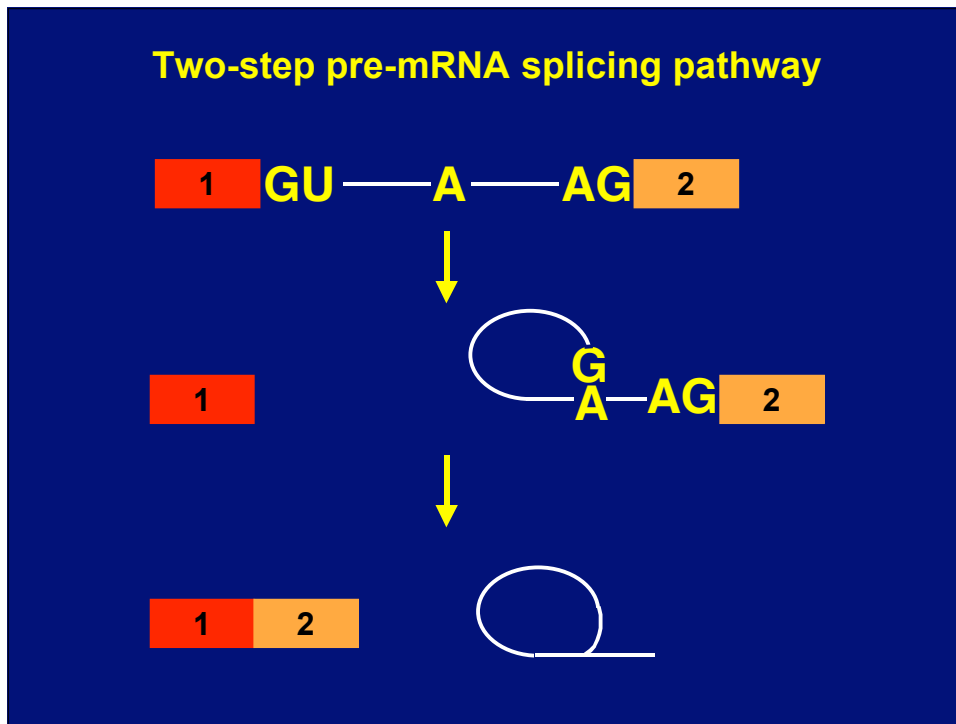## Classical view of eukaryotic gene expression

# An ESE (enhancer-sequence-element) SNP can alter RNA splicing

(Collaborated with Krainer lab at CSHL: *Nature Genet*. 2001)

## Patterns of alternative pre-mRNA splicing

Exon inclusion/skipping

Mutually exclusive exons

Alternative 3' splice sites

Alternative 5' splice sites

Retained intron

Alternative promoters

Alternative polyadenylation sites

## Two-step pre-mRNA splicing pathway

## Conserved intronic splicing signals

5' splice site          branch site          3' splice site

Adapted from Burge et al (1999) in *The RNA World*

**Functional SELEX approach to study SR protein specificity**

SR

IgM

M1 exon    intron    M2 exon    ESE

IgM splicing in S100 extract requires

- SR proteins
- Exonic Splicing Enhancer



**ESE consensus motifs**

(G&D1998,MCB1999)

SF2/ASF    Max : 6.589    Thr: 1.956

SC35    Max : 6.221    Thr: 2.383

SRp40    Max : 6.324    Thr: 2.670

SRp55    Max : 6.135    Thr: 2.676

**A *BRCA1* nonsense mutation causes exon skipping (NMD?)**
(Mazoyer et al Am J Hum Genet 1998)

Mazoyer et al (1998) *Amer J Hum Genet* 62: 713-715



**Mutants to test nonsense codon requirement**

SF2/ASF score

WT  AT GCT GAG TTT GTG  2.143

NL  AT GCT tAG TTT GTG  0.079

NH  AT GCT GAG Tag GTG  2.496

ML  AT GCT aAG TTT GTG  0.079

**Different SF2/ASF high-score motifs can substitute for the natural ESE**



Alteration of enhancer motif scores by point mutations (50 exon-skipping mutants)

ESEfinder: a web resource to identify ESEs
(http://exon.cshl.edu/ESE/ Cartegni et al. submitted)



• *ASDB* (with Stamm/Max-Plank, Nakai/Kyoto, 2001)

• *mATDB* (Wang & Zhang, to be submitted)

• **The RNA-mediated Annealing, Selection and Ligation Assay (RASL): Application in Alternative Splicing** (with Fu&Gribskov/UCSD and Fan/Illumina, NCI funded)

**Step 1**: Oligo Annealing

**Step 2**: Poly A$^+$ Selection

**Step 3**: RNA-dependent Oligo Ligation

**Step 4**: Amplification using Universal Primers

**Step 5**: Hybridization to Zip-Code Arrays

# Classification of 5'UTRs by *CART*

## (Collaborated with Sunoga lab at Tokyo U. *Genome Res*. 2000)

### Forms of translation regulation by 5'UTR

| | | | | | |
|---|---|---|---|---|---|
| m7G | secondary | AUG----UAA | AUG | IRES | translation start |
| tss Site | structure | (uORF) | (uAUG) | | (AUG) |

**Some examples**

| Gene Name | Length | uAUGs | Mechanism |
|---|---|---|---|
| c-mos (ovarian mRNA) | 80 | 0 | secondary structure |
| c-mos (testicular mRNA) | 300 | 4 | uAUG |
| RAR beta2 | 461 | 5 | uORF |
| PDGF2/c-sis | 1022 | 3 | IRES |
| TGF beta 1 | 840 | 0 | secondary structure |
| ATM (cancer) | 146-884 | 1 to 8 | uAUG |
| AR | 1116/1127 | 1 | uAUG |
| c-myc | 408 | 0 | IRES |
| FGF-2 | 484 | 0 | IRES |
| IGFII (p3) | 1170 | 0 | secondary structure |
| IL-15 | 313 | 10 | uAUGs |
| TGF beta 3 | 1104 | 11 | uAUGs |
| TGF beta 3 (Breast cancer cells) | 297 | 0 | highly expressed |
| Spi-1 | 151 | 0 | secondary structure |

# Computational Molecular Biology of Genome Expression and Regulation

**5'UTR database**
- a set of 954 human 5'UTR sequences was obtained from *5' end-enriched cDNA library* (Suzuki et al. 2000) with their mRNA start sites mapped
- a second set of 1613 full-length 5'UTR sequences retrieved from UTRdb (Pesole et al. 2000) database
- all the redundant and ambiguous sequences were eliminated and finally a non-redundant set of 2312 5'UTR sequences was prepared for the analysis

**CART classification of human 5'UTR sequences**
- *Class I(226):* 5'UTRs of growth factors, their receptors, transcription factors, proto-oncogenes, cytokine receptors and tumor suppressor genes. Most of these are understood to be **translationally repressed** mRNAs.

- *Class II(70):* This class consists of TOP mRNAs.(5'terminal oligopyrimidine tract-5'TOP), The **translation is regulated in growth dependent manner**.

- *Class III(76):* 5'UTRs of highly expressed genes, tubulins, globins, globulins, myosins, caseins, glycolytic enzymes, beta-actin, gamma-actin and histones. Theses transcripts are believed to be **efficiently translated or (at least) not repressed at the translational level**.



**CART Model** (Breiman,Friedman, Olshen & Stone, 1983)

**Computational Molecular Biology of Genome Expression and Regulation**

## Cross Validation Classification:

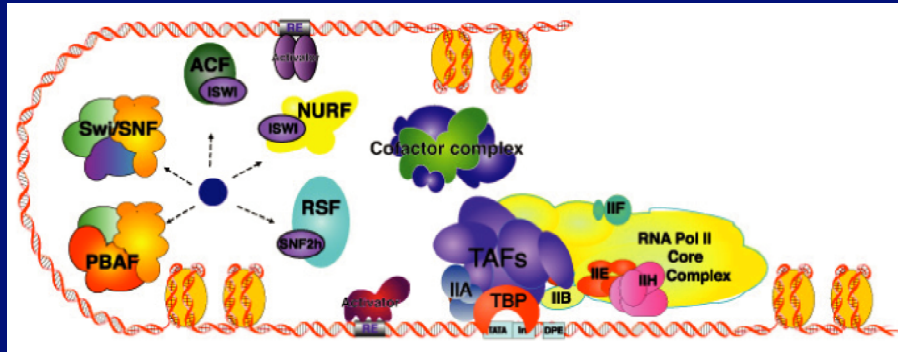| Actual Class | Predicted Class | | | Actual Total |
|---|---|---|---|---|
| | I | II | III | |
| I | 210 (93.0) | 1 (0.4) | 15 (6.6) | 226 |
| II | 0 (0) | 70 (100) | 0 (0) | 70 |
| III | 10 (13.2) | 2 (2.6) | 64 (84.2) | 76 |
| Predicted Total | 73 | 220 | 79 | 372 |

**RESULTS**

**ΔG was the most discriminative variable for the three classes.** ΔG was followed by presence of TOP, 5'UTR length, number of stable free energies, presence of stable secondary structure within the first 100 bp from the cap site, CDS length, A/T ratio, G/C ratio, number of uAUGs, GC%, number of uORFs and codon bias, in the order of relative importance for predictive classification.

**More than 90% of Class I 5'UTRs are embedded with stable secondary structures with ΔG less than –50 kcal/mol.** Classes II & III are almost free from this translational inhibitory feature. Also 60% of the Class I 5'UTRs have stable secondary structures within the proximity of the transcription start site.

**Presence of uAUGs and uORFs was observed as common feature in Class I 5'UTRs** where as Classes II and III are quite free from these features.

**65% of Class II transcripts are in good start site context** followed by Class III with 57% and Class I with 49%.

**There was not any significant difference in GC% between the three classes.**

There was not any significant difference in mean codon bias between the three classes, which indicates that **the codon usage and expression level in human genes are not correlated.** In contrast, codon bias plays an important role in translational efficiency in some lower eukaryotes, such as yeast (Sharp and Li, 1987).

# Promoter analysis *in silico*



- •De novo (database/training set) -> TSS;
- •Functional genomics (expression/localization) -> cis-elements;
- •Comparative genomics -> TSS & cis-elements

# Core_Promoter (Zhang, Genome Res. 1998)



Pol.IID (TFs)HETF150FTF250BXUPEATATATBPCIFDPEInr-60 — -30

GenBank
M12523:1..1980
ALB gene=serum albumin
Firstexon=1737..1854
AUG=1776
C+G=0.33

**Core_Promoter prediction**:

| TSS | Score |
|------|-------|
| **1737** | **0.637** |
| 1736 | 0.604 |
| 1727 | 0.588 |
| 1732 | 0.534 |
| 1731 | 0.531 |
| 1728 | 0.498 |
| 1726 | 0.428 |



QDA variables

## CpG_Promoter (Ioshikhes&Zhang *Nature Genet.* 2000)

CpG island: Length > 200 bp;C + G content > 50%;CpG ratio Obs/Exp > 0.6

- 135 genes
- 68 have CpG island around promoter
- 63 recognized
- SN = 0.47 (0.93)
- SP = 0.34 (1 Pos./26 kb; 1/36 kb is in fact)

•Promoter Scan gives
  SN = 0.44
  SP = 0.06 (1 Pos. / 4.7 kb)

| GenBank | CpG_Promoter prediction: | | Core_Promoter prediction: | |
|---------|--------------------------|---|---------------------------|---|
| | CpG islands associated | Promoter- | TSS | Score |
| D87675 | | | **8921** | **0.100** |
| >301kb | **8813..9319** | + | 8923 | 0.094 |
| App gene encodes | **9328..9547** | + | 8920 | 0.089 |
| Amyloid precursor protein | **9761..10203** | + | 8919 | 0.084 |
| Firstexon=9001..9204 | 117256..117511 | - | 8922 | 0.078 |
| AUG=9148 | 176132..176342 | - | 8918 | 0.058 |
| | 257735..257942 | - | 8783 | 0.056 |
| | 261475..261750 | - | | |

## First exon prediction (FirstEF)
### (Daluvuri,Grosse&Zhang, Nature Genet. 2001)

**Performance statistics of FirstExonFinder based on cross validation**

| Exon Type | Sn | Sp | CC |
|-----------|-----|-----|-----|
| CpG-related | 0.92 | 0.97 | 0.94 |
| Not CpG-related | 0.74 | 0.6 | 0.65 |
| All Exons | 0.86 | 0.83 | 0.83 |

**Promoter Prediction accuracy of FirstEF and PromoterInspector for Ch22**

| Program | TP | FP | Sn | Sp |
|---------|-----|-----|--------|--------|
| FirstEF | 46 | 40 | 79.30% | 53.50% |
| PromoterInspector | 28 | 37 | 48.30% | 43.10% |

**Prediction Accuracy for Ch21&22 (Number of Real Promoters: 58)**

| Chromo-some | Number of Exponentially mapped first exons | Number of correctly predicted first exons | Completely non-coding exons | Predicted non-coding exons |
|-------------|--------------------------------------------|-------------------------------------------|-----------------------------|----------------------------|
| 21 | 42 | 37 (88%) | 14 | 10 (71%) |
| 22 | 79 | 69 (87%) | 28 | 23 (82%) |
| Total | 121 | 106 (88%) | 42 | 33 (79%) |

## FirstEF program



http://www.cshl.org/mzhanglab/FirstEF

## Cell Cycle Regulation (continued)

**Computationally predicted E2F target genes confirmed by**
*in vivo* **footprint (ChIP) (Kel** *et al.* **JMB 2001)**

| Gene | EMBL | Sequence of the potential sites | Position rel. start transcription | Score, q | d(X) |
|---|---|---|---|---|---|
| **c-fos**, *Homo sapiens* | HSFOS | (-) gtCTTGGCGCGTGTcc<br>(-) ggGGTGGCGCGCGGgc<br>(+) ccTCTGGCGCCACCgt<br>(-) acGGTGGCGCCAGAgg | -165 .. -176<br>-92 .. -103<br>-90 .. -79<br>-78 .. -89 | 0.915<br>0.836<br>0.878<br>0.830 | 2.92 |
| **JunB**, *Homo sapiens* | HS207341 | (+) gcGATCGCGCCAGAga<br>(+) tcTCTGGCGCGATAgc<br>(-) ggGCTGGCGCGGGCgg | 79 .. 90<br>91 .. 80<br>169 .. 158 | 0.887<br>0.905<br>0.820 | 3.16 |
| **TGF-β1**, *Homo sapiens* | HS1GFB1PR | (+) ctGGTTGGCGGGGGCGga<br>(+) ccCTTCGCGCCCTGgg<br>(+) ctCTTGGCGCGACGct<br>(-) agCGTCGCGCCAAGag<br>(+) ccTTTGCCGCGGGGga | -513 .. -502<br>-298 .. -287<br>28 .. 39<br>40 .. 29<br>85 .. 96 | 0.804<br>0.912<br>0.928<br>0.830<br>0.854 | 2.03 |
| **ARF**, *Homo sapiens* | AF082338 | (+) acTTTGCCGGCCCTGtg | -265 .. -276 | 0.859 | |
| **Mcm4** (Cdc21), *Mus musculus* | AB000629 | (+) ggTTTCCCGCCAAAAcc<br>(-) gtTTTGGCGGGAAAcc<br>(+) gcAGTGGCGCCTTCcg<br>(-) ccTTTGCCGCTGTGat<br>(-) tgGGTGGCGCAGAAct<br>(+) ttTGTCGCGCAGCAac | -443 .. -432<br>-431 .. -442<br>-329 .. -318<br>-297 .. -286<br>-127 .. -116<br>-24 .. -13 | 0.872<br>0.935<br>0.810<br>0.846<br>0.809<br>0.858 | |
| **MCM5** (P1-CDC46), *Homo sapiens* | HS286B10 | (+) agTTTCGCGCCAAAAct<br>(-) aaTTTGGCGCGAAAct<br>(+) ttTTTCCCGCGAAAct<br>(-) agTTTCGCGGGAAAaa | -187 .. -176<br>-175 .. -186<br>8 .. 19<br>20 .. 9 | 0.988<br>1.005<br>0.885<br>0.932 | 4.91<br><br>3.01<br>4.21 |
| **von Hippel-Lindau** (VHL), *Homo sapiens* | AF010238 | (+) aaGCTCGCGCCACTgc<br>(-) gcAGTGGCGCGAGCtt<br>(-) gtCTTCGCGCGCGCtc | -270 .. -259<br>-258 .. -269<br>-28 .. -39 | 0.810<br>0.838<br>0.921 | 2.22 |
| **B-myb**, *Homo sapiens* | HSBMYBDNA | (-) gtCCTGGCGCGCGGgc | -72 .. -83<br>-53 .. -42 | 0.831<br>0.866 | 5.50 |
| **Nucleolin**, *Homo sapiens* | HSNUCLEO | (-) ttTTTGGCGCCGGCtg<br>(-) ccGTGGGCGCGCGGgt | -297 .. -308<br>-256 .. -267 | 0.966<br>0.814 | 2.91 |
| **Nucleolin**, *Cricetulus griseus* | CSNUCLEO | (-) cgTTTGGCGCGGCTtg | -296 .. -307 | 0.973 | 6.67 |
| **Nucleolin**, *Mus musculus* | MMNUCLEO | (-) agTTTGGCGCGGCTtg | -306 .. -317 | 0.973 | 1.76 |

---

## Does genome location analysis work with human cells?
—— Challenges (In collaboration with Ren Lab UCSD)

**ChIP-chip**:

- Genome is 200 times bigger than yeast

- Abundant repetitive sequences

- Annotation of gene structure and function is much less complete

- Many different cell types

- Quality of antibodies

## Location Analysis in Human cells:
### E2F4 Binding to Human Promoters (Ren et al. 2002)

$P = 0.01$

Total Input DNA
Cy3 (Total DNA)

E2F4 ChIP DNA
Cy5 (ChIP DNA)

E2F3
MCM3
CDC6
p107
CDC25A
ORC1
CYCLIN A2
CDC2
POL A
E2F2
Rb

## RT-PRC Confirmation (Cont.)

**DNA Repair**

| | Wt | p107[-/-] p130[-/-] |
|---|---|---|
| Hr | 0 8 12 20 | 0 8 12 20 |

CYCLIN A
ACTIN

BARD1
RAD54L

**Checkpoints**

| | Wt | p107[-/-] p130[-/-] |
|---|---|---|
| Hr | 0 8 12 20 | 0 8 12 20 |

MAD2L
TTK
CHEK1

**Mitotic**

| | Wt | p107[-/-] p130[-/-] |
|---|---|---|
| Hr | 0 8 12 20 | 0 8 12 20 |

HEC
NEK2
SECURIN/PTTG1

New challenge: Expression+ChIP+Promter, CompositeSites

Comparative DNA sequence analysis of protocadherin gene clusters
(Maniatis, Meyers, Zhang labs, Genome Res. 2001)

**Computational Molecular Biology of Genome Expression and Regulation**



Mouse and Human CNS Clusters



Phylogenetic relationship

**Upstream Sequence Similarity**



**CpG-islands, conserved promoter elements**



Pipmaker                    Gibbs sampler

## Novel conserved regions in the introns        (transcplicing signals?)

**A**        Conserved sequences upstream of Pcdhα constant region exon 1 (83%)

```
M: -949 tagccaaacaaggtcaaacta--cagtggtttgttttcttctctttgctgtccctcctcc
        |||| |||||||||||||| |||||| | ||||||||| |||| | | ||||||| ||
H: -604 tagcaaaacaaggtcaaactctgcaatagtttgttttcctctccctagtatccctcttca

        atcagttgaagaggctgttgtcagttgctagtgttacgactgggcacatccttctcaggt
        ||||||  ||||| ||||| ||||||||| |||||| |||||||||||||  |   |||
        atcagaa-aagagactgttatcagttgctggtgttatgactgggcacatccgccctgggt

        caaacatgctgcagtctgcaaagccagcagtagattgcagtcctctgcagtccaggcaga
        |||| ||||||||||||||||||||||||| ||||||||||||||||||||||||| |||
        caaatatgctgcagtctgcaaagccagcagcagattgcagtcctctgcagtccagccagg

        tctgcagaatttgtgt  -756      **B**      Conserved sequences upstream of Pcdhγ constant region exon 1 (83%)
         | |||||| |||||
        ccagcagaacttgtgt  -410    M: –671  tttctccgccttggacagagctgccttgttcccatctccttggtcacaggccattgtgag
                                            |||||||| || ||||||||| |||||| ||| ||||| |||||| ||||||| |
                                  H: –714  tttctccgtctcagacagagcagccttgttctcttctccttagtcacagaccattgtctg

                                           gcatgaagttctgggggtgagaggcatcccagagcttggatgccctataaaggctcaggc
                                           ||| | |||||| ||||||||| | |||| | |||||||||| ||||||| | |
                                           gcacggagttctaggggtgagaagtgtcccgggacttggatgccccgcaaaggcccaatc

                                           tggcatgactcctaaattaataatgtatttagctgtgggaagag-gtctttgagatcgag
                                           |||||||||||||||||||||||||||||||||||||||||||| |||| |  | |||
                                           tggcatgactcctaaattaataatgtatttagctgtgggaagagattcttgcaagccaag

                                           ggcccggaggaggtggccctctgaatgtgtgagtgcacaacgtggcacaaaaagggttac
                                           ||||| ||| || || |||| ||||||||| ||| |||||||| |||||| |||||||||
                                           ggcccagagaagatgtccctgtgaatgtgtcactgcacaacctggcaccaaaagggttac

                                           ccagagcagcagccatcttgctgcagccgaggctttgttcccagctgaggagctgaat -375
                                           | ||| |||||||||||||||||||| || |||||||||||||||||||||||| ||||
                                           caagaacagcagccatcttgctgcagaggatgctttgttcccagctgaggagttgaat -417
```

## Human genome distribution of CSEs

Xuan, Wang & Zhang, Genome Biology (2002)

CSEs and Human Gene number *(cont.)*



Zhang Lab Members