

## Supporting Methods

### Sampling error

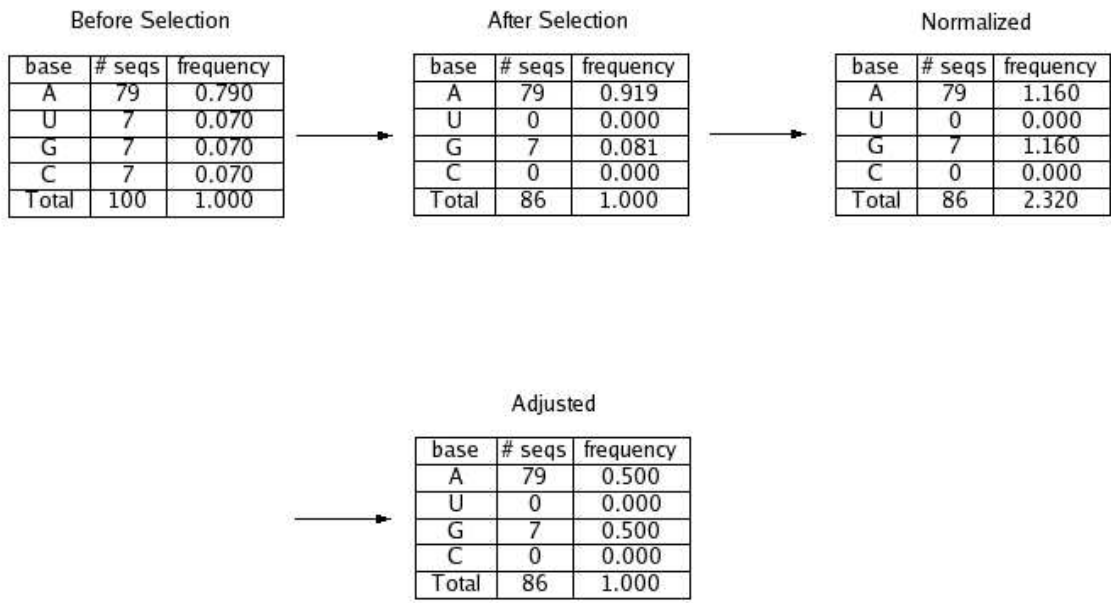
Treating the per-position nucleotide frequencies ( $F_i$ ) as though they were actual probabilities ( $P_i$ ) results in a sampling error in the information content calculation. Schneider et al.<sup>21</sup> demonstrated a method for calculating the standard error by computing the weighted uncertainties for all possible combinations of bases for a given number of sequences in an alignment. To account for the reduction in uncertainty due to a finite number of sequences, we subtract a correction factor from the information content calculated at each position. The correction factor is the difference between the known maximum uncertainty (2 bits/position for RNA) and the calculated expectation value for the maximum uncertainty in an alignment of a certain size. To compute the values for the expected maximum uncertainty, error correction and variance per position, we implemented a version of CALHNB by Schneider et al.<sup>21</sup> in Perl. We modified their algorithm to adjust the frequency distributions obtained from the sequence alignments to compute the information content required to specify the active structures in completely random RNA sequence space. The values were computed for all of the alignments with prior base probabilities of 0.79, 0.07, 0.07, 0.07. The chart below gives a sense of the magnitude of these values.

<u>Number of sequences</u>	<u>Expectation value (bits)</u>	<u>Error correction (bits)</u>	<u>Variance per position</u>
30	1.74	0.26	0.060
40	1.82	0.18	0.034
50	1.86	0.14	0.020
60	1.89	0.11	0.013
70	1.91	0.09	0.009
80	1.92	0.08	0.006

## Normalizing and adjusting base frequencies

We computed the information content of each aptamer using the selection sequence alignments. The sequence pools were synthesized with 21% rate of mutagenesis, where each of the non-wild-type bases was equally likely. Because we know the prior probability of each base type  $i$  at every position, we can use the frequency distribution of nucleotides from the alignment to compute the information content of the aptamers in relation to random ( $P_A=P_U=P_C=P_G$ ) sequence space.

As an illustration (see below), consider an aptamer that requires a purine (A or G) in a certain position to be functional. For simplicity imagine that we synthesize a pool of 100 aptamer sequences mutated 21% at that position with A as the original base. In the idealized case we would have 79 sequences with an A, 7 with a U, 7 with a G, and 7 with a C. After selection, only the A- and G-containing sequences remain. The total number of sequences in the pool is reduced to 86, where 79 are A and 7 are G. The relative frequency of A and G after the selection reflects the fact that A and G are equally compatible with a functional molecule and that A was present much more often in the starting pool than G. To normalize each  $F_i$  for the skewed starting frequencies we divide each  $F_i$  after the selection by the  $F_i$  before the selection. Dividing each normalized  $F_i$  by the sum of the normalized  $F_i$ s gives the adjusted frequencies of each base type relative to the others. In the current example, the normalized and adjusted  $F_A=0.5$  and  $F_G=0.5$ ; the information content of this position is 1 bit, which is the information content needed to specify a two-base varying position in an RNA structure.



**Supporting Figure 1A-K.** GTP aptamer sequence alignments. DNA sequence alignments from each of the selections are shown in order from highest affinity for GTP to lowest. The top line in each alignment contains the minimized sequence (corresponding to the length of the optimized sequence). This is the sequence used as the template for synthesizing the library of mutants. Immediately below that is the secondary structure model of the optimized aptamer shown in bracket form where the 5' -most open bracket "(" base pairs with the 3' -most closed bracket ")". Paired strands are shaded with same color. Point mutations in the loops that were assayed are indicated on the same line as the secondary structure. Blue bases were tested and found not to improve binding activity. Red bases improved the activity of the aptamer relative to the minimized construct and were incorporated into the optimized version.

Mutations in the stem regions of the alignment are shaded to indicate:

Red = Watson-Crick covariation  
Orange = new Watson-Crick base pair  
Green = new wobble (A-C or G-U) pair  
Black = broken pairing

Mutations in the loop are shaded purple





Figure 1C. 10 – 10.

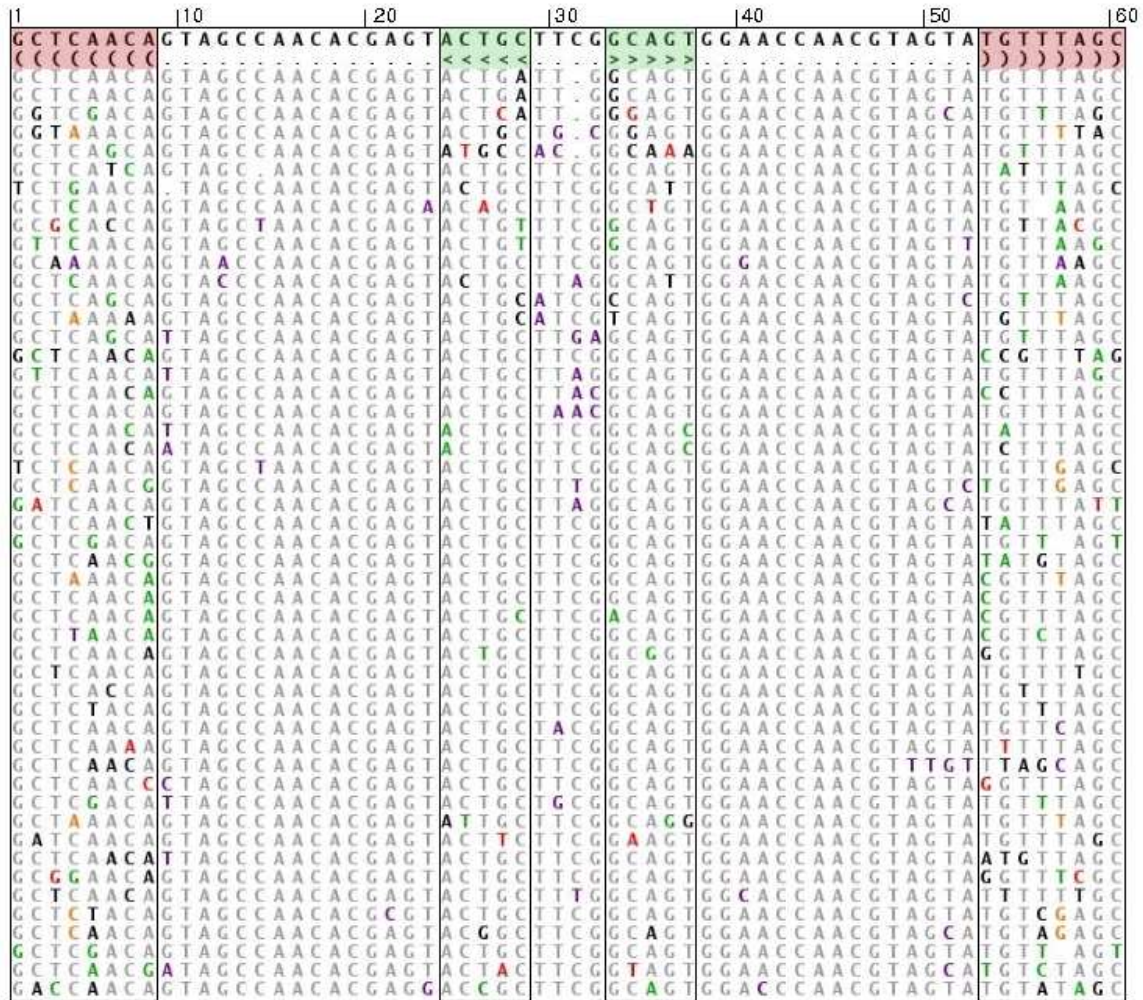










Figure 1G. 9 – 12.

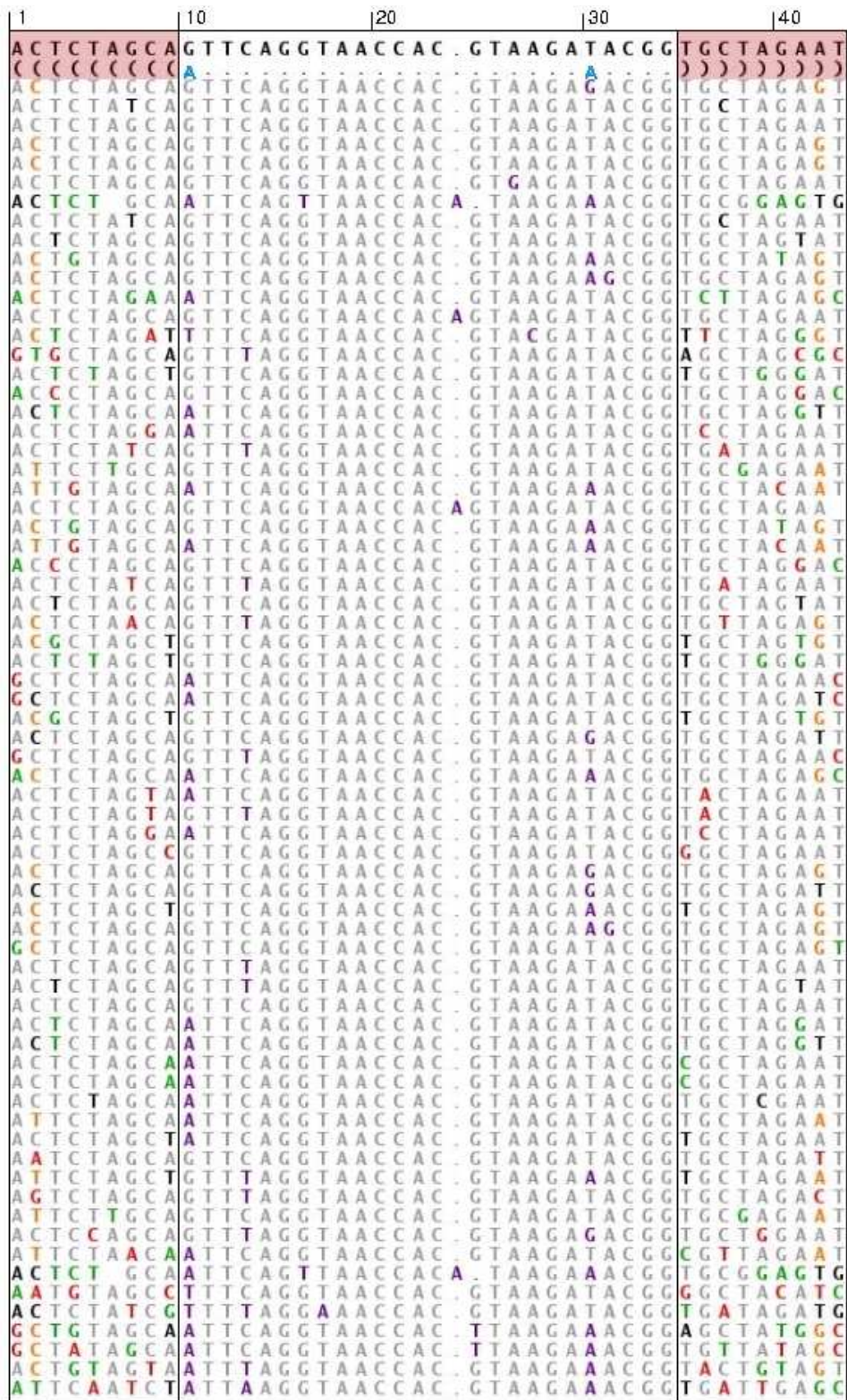


Figure 1H. 10 – 6.

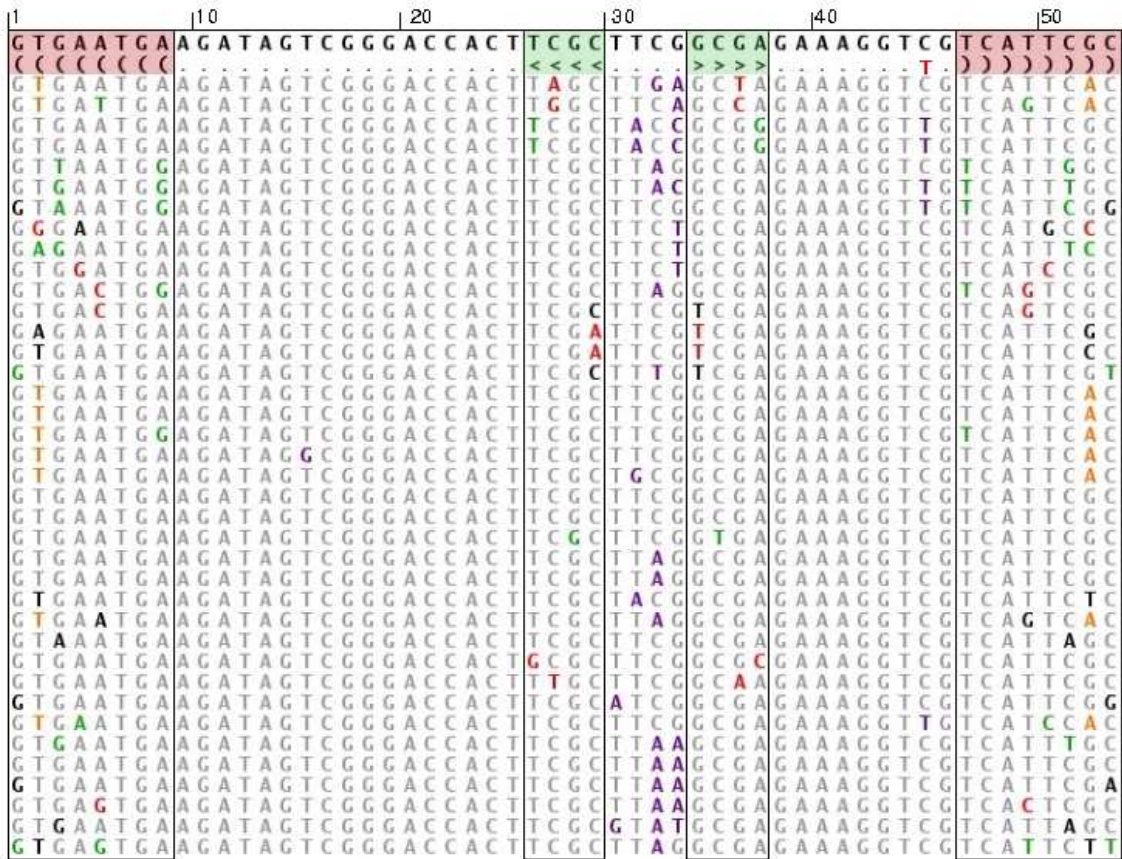




Figure 1J. Class IV.

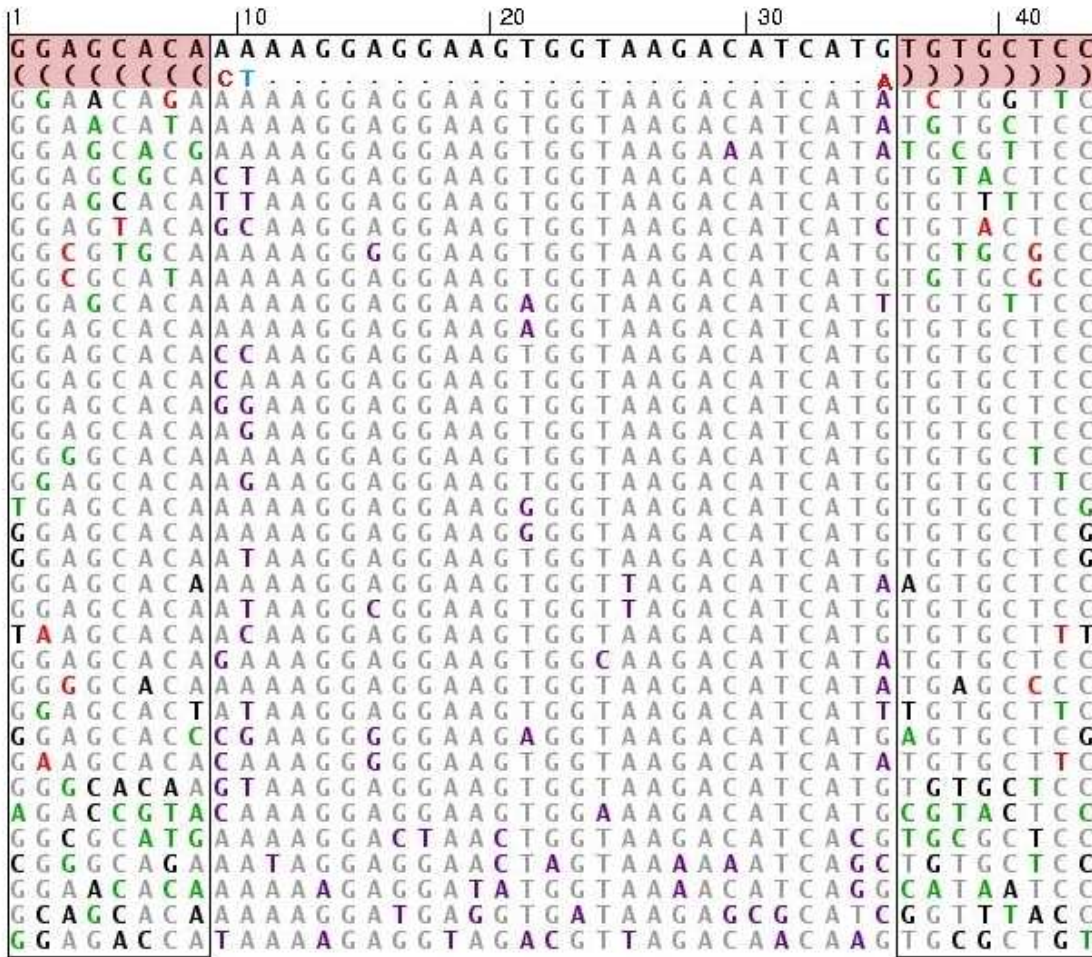


Figure 1K. Class III.

1	10	20	30	40
C C C A T G A A	C G T C	T T C G	C A G T	T T G C T A A A A A C C A G T C . G T G G G
( ( ( ( ( ( ( (	< < < <	- - - -	> > >	) ) ) ) )
C C C A T T T T	T G T G C	T T C G	C A G T	T T T C T A A A A T A C A A G T C . G T A G T
C T C A T G A A	C G G T	T T T G	C C C G	T T G C T A A A A A C C A G T . G T G G G
C C C A T G A A	C T G A	T T T G	T C C G	T T G C T A G A A A A C C A G T . G T G G G
C C A T T G C A	C G G A	T C C G	C C C G	T T G C T A A A A A C C A G T . G T G G A
C T C A T G A T	C G T C	T T C G	C A C G	T T G C T A A A A A C C A G T C . G T G A G
C T C A T G A T	C G T C	A T C G	T A C G	T T G C T A A A A A C C A G T C . G T G A G
T C C A T G A T	C G T C	C T C G	C A C G	T T G C T A A A A A C C A G T C . G T G G A
T C C A T G A T	C G T C	C T C G	C A C G	T T G C T A A A A A C C A G T C . G T G G A
T C C A T G A T	C G T C	C T C G	C A C G	T T G C T A A A A A C C A G T C . G T G G A
T C C A C G A T	C G T C	C C C G	C A C G	T T G C T A A A A A C C A G T C . G T G G G
C C C A T G A T	C G T C	T G C G	T A C G	T T G C T A G A A A A C C A G T C . G T G G G
C C C A T G A T	C G C C	T T C G	T A C G	T T G C T A A A A A C C A G T C . G T G G G
C A C A T G A T	C G T C	T T C G	C A C G	T T G C T A A A A A C C A G T C . G T G T G
C C C G T G A T	C G C C	T T C G	C A C G	T T G C T A A A A A C C A G T C . G T G G G
C C C A T G A T	C G A C	T T C G	A T C G	T T G C T A A A A A C C A G T C . G T G T G
C C A T T G A T	C G T C	A T C G	C A C G	T T G T T A A A A A C C A G T C . A T C G G
G C C A T G A T	C G T C	T T C G	A A C G	T T G C T A A A A A C C A G T T . G T G G C
C T T A G G A T	C A G C	T T T G	C C T G	T T G C T A A A A A C C A G T C . C T G G C
C C A A G G A T	C A G C	T T C G	C C T G	T T G C T A A A A A C C A G T C . C T G G T
C C C A T G A T	C A T C	T T C G	C A T G	T T G C T A A A A A C C A G T C . C C G G G
C C C A T G A T	C A T C	T T C G	C A T G	T T G C T A A A A A C C A G T C . C C G G G
C C C C T G A T	C A T C	T T T G	A A T G	T T G C T A A A A A C C A G T C . G A G G G
C C C A T G A T	C A G C	T T C G	C T G	T T G C T A A A A A C C A G T C . A T G G G
C C C T T G A T	C A G C	T T C G	C T G	T T G C T A A A A A C C A G T C . G T G G G
C C C A T G A T	C A G C	T T C G	C C T G	T T G C T A A A A A C C A G T C . G T G G G
C C C A T G A T	C A G C	T T C G	C C T G	T T G C T A A A A A C C A G T C . G T G G G
C C C A A G A T	C A G C	T T C G	C C T G	T T G C T A A A A A C C A G T C . G T G G G
C T C A T G A T	C A G C	C T C G	C C T G	T T G C T A A A A A C C A G T C . G T G A G
C C C A T G A A	C T G C	T T C G	G C A G	T T G C T A A A A A C C A C T C . G T G G G
C C C T T G A T	C A G C	T T T G	C T G	T T G C T T A A A A C C A G T C . T T G G G
C C C A A G A T	C A G C	G T C G	T C T G	T T G C T A A A A A C C A G T C . T T G G G
C C C A T T A T	C A G C	T T C G	C C C G	T T G C T A G A A A A C C A G T G . T A G A G
C C C A T G A T	C A G C	C A T C	G C T G	T T G C T A A A A A C C A G T C . A T G A G
A C C A T G A A	C T C C	T C C G	C C C G	T T G C T A A A A A C C A A T T . G T G G G
C T C A C G A T	C A T C	T T C G	A A T G	T T G C T A A A A A C C A G T C T G T T T C
C A C A T G A T	C G T C	C T T G	C A C G	T T G C T A A A A A C C A G T C T G T G T A
C C C A T G A T	C G T C	A T C G	C A C G	T T G C T A A A A A C C A G T C T G T G G

**Supporting Chart 1.** Original, minimized, and optimized sequences. The chart contains DNA sequences for each aptamer included in our analysis. The top in each group, marked 'ori', is a clone that was identified from the original GTP aptamer selection of Davis and Szostak<sup>19</sup>. Aptamers referred to by number indicate that only one isolate was identified. Aptamers named Class I-V originated multiple times as judged by the presence of diverse flanking sequences - these suggested 5' and 3' boundaries for the functional regions ('ori K<sub>d</sub> construct'). For aptamers with multiple isolates (Class I-V), the specific original sequence chosen to derive the optimized versions is displayed. The 'min K<sub>d</sub> construct' is the minimized functional sequence determined by deletion end mapping. 'min' is the minimized sequence flanked by constant primer regions and was employed as the template for selecting active sequence variants. The last sequence, 'opt', is the optimized version of each aptamer. K<sub>d</sub>s are shown in parentheses. The constant primer-binding sites included for RT-PCR are shown in blue, where applicable. Sequence changes, relative to the original isolate, are shown in red. Occasionally unpaired G's were appended to the 5' ends of the 'min K<sub>d</sub> construct' and 'opt' constructs to increase transcription yields; these are colored green.

The optimized sequences will be deposited in GenBank. All sequences, (including selection alignments) are available in FASTA format on our website:

<http://genetics.mgh.harvard.edu/szostakweb/publications/extra/informationalcomplexitypaper2003.html>



## 9-4

ori (25 nM)

5' -GGAGGCGCCA**ACTGAATGAA**AGTTGCCAGCTGCGAGCACGTGAATAGACTGCTTCGGCAGTGTCTCG  
ACGTGTGTAGGGGAAAGTATCCT**CCGTA**ACTAGT**CGCGTCAC**-3'

min (K<sub>d</sub> construct 25 nM)

5' -GGAGGCGCCA**ACTGAATGAA**AGTTGCCAGCTGCGAGCACGTGAATAGACTGCTTCGGCAGTGTCTCG  
ACGTGTGTAGGGGAAAGTATCCT**CCGTA**ACTAGT**CGCGTCAC**-3'

min

5' -GGAGC**ACGAACTCGGTATCC**GGAGGCGCCA**ACTGAATGAA**AGTTGCCAGCTGCGAGCACGTGAATAG  
ACTGCTTCGGCAGTGTCTCGACGTGTGTAGGGGAAAGTATCCT**CCGTA**ACTAGT**CGCGTCACGCGAACC**  
**ATTCGGAACATCG**-3'

opt (9 nM)

5'-GG**GACG**AGCACGTGAAT**CGACTGCTTCGGCAGTGTCTCGACGTGTGTAGGGGAAAGTATCC****CCCGTCC**  
**C**-3'

## Class V

ori (4000 nM)

5' -GGAGGCGCCA**ACTGAATGAA**ATTTGGGCATTTTGGTAGGTCGGTCGCTGCTTCGGCAGTAAGGGGT  
AGGCATTGCTGGCCTAGGGT**CCGTA**ACTAGT**CGCGTCAC**-3'

min (K<sub>d</sub> construct 15000 nM)

5' -GGG**GGGCATTTTGGTAGGTCGGTCGCTGCTTCGGCAGTAAGGGGTAGGCATTGCTGGCCTAGGGTC**  
CGTAAC-3'

min

5' -GGT**ATACGTGCAGAGACGCG**TTGGGCATTTTGGTAGGTCGGTCGCTGCTTCGGCAGTAAGGGGTAG  
GCATTGCTGGCCTAGGGTCCGTAAC**GCATGATAGCTGATCGCAGC**-3'

opt (17 nM)

5' -GGGGGCATTTTGGTAGGTCGGTCGCTGCTTCGGCAGT**G**AGGGGTAGGCATTGCTGGCCTAGGGTCC  
**CC**-3'

## 10-10

ori (30 nM)

5'-GGAGGCGCCA**ACTGAATGAA**TTGCTCAACAGTAGCCAACACGAGTACTGCTTCGGCAGTGGAACCAAC  
GTAGTATGTTTAGCAT**TCCGTA**ACTAGT**CGCGTCAC**-3'

min (K<sub>d</sub> construct 30 nM)

5'-GG**TTGCTCAACAGTAGCCAACACGAGTACTGCTTCGGCAGTGGAACCAACGTAGTATGTTTAGCAT**-  
3'

min

5'-GGAGC**ACGAACTCGGTATCC**TTGCTCAACAGTAGCCAACACGAGTACTGCTTCGGCAGTGGAACCAAC  
GTAGTATGTTTAGCAT**GCGAACCATTTCGGAACATCG**-3'

opt (30 nM)

5'-GG**GCTCAACAGTAGCCAACACGAGTACTGCTTCGGCAGTGGAACCAACGTAGTATGTT****GAGC**-3'

## Class I

ori

5'-GGAGGCGCCAAC**TGAATGAA**UUGCUUCGAGTCTTGAAGTGGTTGGGCTGCTTCGGCAGTGTGAAAATG  
AGGCTTTTAAAGGG**TCCGTA**ACTAG**TCGCGT**CAC-3'

ori (  $K_d$  construct 160 nM)

5'-GGGAGTCTTGAAGTGGTTGGGCTGCTTCGGCAGTGTGAAAATGAGGCTTTT-3'

min ( $K_d$  construct 76 nM)

5'-GGGA**CG**AAGTGGTTGGGCGCTTCGGCGTGTGAAAA**CGTCTC**-3'

min

5'-GGAGCACGAACTCGGTATCCGGGA**CG**AAGTGGTTGGGCGCTTCGGCGTGTGAAAA**CGTCTC**GCGAAC**C**  
ATTCGGAACATCG-3'

opt (76 nM)

5'-GGGA**CG**AAGTGGTTGGGCGCTTCGGCGTGTGAAAA**CGTCTC**-3'

## 10-59

ori (285 nM)

5'-GGAGGCGCCAAC**TGAATGAA**CAGGATGGTAAGTTCCCAAGGCGGGTTGGAAGAGATATCATAGGAGCT  
TGTCGTTCTGGTCATCCT**TCCGTA**ACTAG**TCGCGT**CAC-3'

min (  $K_d$  construct 285 nM)

5'-GGATGGTAAGTTCCCAAGGCGGGTTGGAAGAGATATCATAGGAGCTTGTCGT-3'

min

5'-GGAGCACGAACTCGGTATCCATGGTAAGTTCCCAAGGCGGGTTGGAAGAGATATCATAGGAGCTTGTC  
GTGCGAAC**CA**TTTCGGAACATCG-3'

opt (250 nM)

5'-GGGATGGTAAGTTCCCAAGGCGGGTTGGAAGAGATATCATAGGAGCTTGTCGT**CCC**-3'

## 10-24

ori (500 nM)

5'-GGAGGCGCCAAC**TGAATGAA**TGATCCTTTGGAGGGGCATCTATATACTGCTTCGGCAGGGAACTCTAC  
TAAGCACCGATGTCAC**TCCGTA**ACTAG**TCGCGT**CAC-3'

min ( $K_d$  construct 300 nM)

5'-GCGGGCATCTATATACTGCTTCGGCAGGGAA**C**CTCTACTAAGCACCGATGTC**CGC**-3'

min

5'-GGGCATGTCCTTTTT**CCTAAT**GCGGGCATCTATATACTGCTTCGGCAGGGAA**C**CTCTACTAAGCACCG  
ATGTC**CGC**GTAATAGCAGCATACTCGGA-3'

opt (300 nM)

5'-GCGGGCATCTATATACTGCTTCGGCAGGGAA**C**CTCTA**A**TAAGCACCGATGTC**CGC**-3'

## 9-12

ori (300 nM)

5'-GGAGGCGCCAAC**TGAATGAA**CAAAGTTGTTAACACCCACTCTAGCAGTTCAGGTAACCACGTAAGATACGGTGCTAGAATGGGAT**TCCGTA**ACTAGT**CGCGT**CAC-3'

min ( $K_d$  construct 300 nM)

5'-GGGACTCTAGCAGTTCAGGTAACCACGTAAGATACGGTGCTAGAAT-3'

min

5'-GGAGCACGAACT**CGGTATCC**ACTCTAGCAGTTCAGGTAACCACGTAAGATACGGTGCTAGAAT**GCGAA**CCATT**CGGA**ACATCG-3'

opt (300 nM)

5'-GGGACTCTAGCAGTTCAGGTAACCACGTAAGATACGGTGCTAGAG**G**T-3'

## 10-6

ori (500 nM)

5'-GGAGGCGCCAAC**TGAATGAA**GATAGTCGGGACCACTTCCCTGAGTTCTGCTTCGGCAGATAACTCCAA GGAGAAAGGTCGTCATT**TCCGTA**ACTAGT**CGCGT**CAC-3'

min ( $K_d$  construct 500 nM)

5'-CGGAATGAAGATAGTCGGGACCACTTCCCTGAGTTCTGCTTCGGCAGATAACTCCAAGGAGAAAGGTC GTCATT**CGC**-3'

min

5'-GGAGCACGAACT**CGGTATCC**CGGAATGAAGATAGTCGGGACCACTTCCCTGAGTTCTGCTTCGGCAGATAACTCCAAGGAGAAAGGTCGTCATT**CGCGCA**ACCATT**CGGA**ACATCG-3'

opt (300 nM)

5'-CGGAATGAAGATAGTCGGGACCACTT**GCTTCGGC**GAGAAAGGTCGTCATT**CGC**-3'

## Class II

ori

5'-GGAGGCGCCAAC**TGAATGAA**GCTGCTAGCGCGGTAGAAAACCGAGCCGGAAGAGCACGTATACGCAGG GCTCAACTACATT**TCCGTA**ACTAGT**CGCGT**CAC-3'

ori ( $K_d$  construct 4600 nM)

5'-GGGACGGCTGCTAGCGCGGTAGAAAACCGAGCCGGAAGAGCACGTATACGCAGGGCTCAACTACA-3'

min ( $K_d$  construct 2500 nM)

5'-GGGACG**CTG**AACGAGCCGGAAGAGCACGTATACGCAGGGCTCAACTAG-3'

min

5'-GGAGCCGACTAATGATTAAT**CTG**AACGAGCCGGAAGAGCACGTATACGCAGGGCTCAACTAG**GCGAGAG** TCCGTAACATCGC-3'

opt (400 nM)

5'-GGGAGCC**A**GGAAGAGCACGTATACGCA**A**GGCTC-3'

#### **Class IV**

ori (2400 nM)

5'-GGAGGCGCCAAC**TGAATGAA**ATTGACTGTTTGAGCATAAAAAGGAGGAAGTGGTAAGACATCATGTAT  
GTTACGTAACGGTTGTCCG**TAACTAGTCGCGTCAC**-3'

min ( $K_d$  construct 2400 nM)

5'-GGGTTGGAGCA**CAAAAAGGAGGAAGTGGTAAGACATCATGTGTGCTCCC**GT-3'

min

5'-GGACGAGA**ACTCGCATAAGC**GTTGGAGCA**CAAAAAGGAGGAAGTGGTAAGACATCATGTGTGCTCCCG**  
TGCATGATAGCTGATCGCAGC-3'

opt (900 nM)

5'- GGGAGCA**CAAAAAGGAGGAAGTGGTAAGACATCATATGTGCTCC**-3'

#### **Class III**

ori

5'-GGAGGCGCCAAC**TGAATGAA**CACACNCCTAAAAGGATACCCATGAACTGCTTCNGCAGTTTGCTAAA  
AACCACTCGTGGGTACNTCCG**TAACTAAGTCGCGTCAC**-3'

ori ( $K_d$  construct 112000 nM)

5'-GGGACACACCCCTAAAAGGATACCCATGAACTGCTTCGCAGTTTGCTAAAAACCACTCGTGGGTACCT  
T-3'

min ( $K_d$  construct 112000 nM)

5'-GGGACACACCCCTAAAAGGATACCCATGAACTGCTTCGCAGTTTGCTAAAAACCACTCGTGGGTACCT  
T-3'

min

5'-GGAGCCG**ACTAATGATTAAT**ACACACCCTAAAAGGATACCCATGAACTGCTTCGCAGTTTGCTAAAA  
CCACTCGTGGGTACCTTCCG**TCCTACATCGGGCATT**C-3'

opt (8000 nM)

5'- GGGATGAT**CGT**CTTCGGAC**CG**TTGCTAAAAACCA**GTCATCCC**-3'