## Supplementary File for Figure 1

**Supplementary Material to *in silico* screen: The *in silico screen* that lead to NFAT5a and alternative potentially lipid-modified transcription factors**

**Supplementary Table 1 and Supplementary Figure 1.1: Taxonomic distribution of NFAT5 isoforms and their sequence conservation in the N-terminal region of human NFAT5a**

**Supplementary Figure 1.2: Sequences of all five human NFAT5a mRNA transcripts in the regions used for qPCR**

**Supplementary Material to *in silico* screen: The *in silico screen* that lead to NFAT5a and alternative potentially lipid-modified transcription factors**

**The *in-silico* screen**

Protein sequence hits from Hidden Markov Model (HMM) searches [1] with protein domain models taken from PFAM [2] and SMART [3] annotated as nucleic acid binding domains against the non-redundant protein sequence database were scored for the presence of an N-terminal myristoylation signal [4,5], a prenylation site [6,7] or the GPI lipid anchor attachment motif [8]. The resulting list of a few thousand hits was to be further

reduced with the goal to come up with singular targets for an experimental follow up. By excluding entries with doubtful sequence termini, we wished to exclude cases of technically correct prediction but possibly non-existing proteins. Considerations of accessibility of cDNA and cell lines restricted the search to protein entries from well-studied model organisms. The requirement for conservation of the predicted lipid anchor sites in other species enhances the reliability of prediction and the possibility to find biological results of more general value.

Finally, the restriction to experimentally verified transcription factor (TF) domains ensures that the candidate under investigation has truly the potential to act as TF. NFAT5a was finally selected as candidate protein for studying the feasibility of nuclear import of lipid-modified TFs because there is experimental evidence for the existence of this isoform [9-12]. Equally important, knowing NFAT5 as the osmotic stress response TF provided us initial hope that testing isoform a for nuclear import upon salt stress is an accessible functional assay [13,14].

**About other TFs with predicted lipid anchors**

The most interesting question for general biology is whether lipid-modified TFs that are imported into the nucleus in a regulated manner are a more general phenomenon and whether reversible palmitoylation as a mobilization mechanism might be found in contexts other than NFAT5a.

It appears impossible to get an exact overview with regard to lipid modified TFs from a renewed *in silico* screen. Despite more than a decade of availability of "complete genomes", new releases of genome assemblies remain accompanied with addition/losses of thousands of proteins, not to speak about isoforms. Nevertheless, several more years of sequencing have improved the data situation since 2004 when we first aggressively unselected hits to come up with the single, hopefully easy-to-test example NFAT5a. Frequently, termini of proteins are described more reliably. The likelihood to find homologues in genomes of alternative species has increased and this improves chances to test the conservation of posttranslational modification sites. Therefore, cases that we *a priori* omitted previously might become worthwhile to be investigated now. For example, the human protein Q8WYA1 (named MOP9, CLIF or BMAL2), a transcriptional regulator with functions in circadian and hypoxia pathways [15], was considered as a potential candidate alongside with NFAT5a. Only two of its isoforms were known in 2004 but eight isoforms are described in the database today. Isoforms 6, 7 and 8 are predicted targets of myristoylation; yet, they do not have an obvious acylation site.

The following five examples all have at least one isoform reported to have an N-terminal glycine after the leading methionine as well as sequentially close cysteines. The proteins are predicted targets for myristoylation at the N-terminal glycines (with NMT/MyrPS [4]) and the cysteines in their vicinity are likely palmitoylation sites [16]. The sequence features are conserved in neighbouring species. (1) The human protein BTBD7 (Q9P203) and (2) the protein BAC20790 from *Oryza sativa* are predicted to possess a BTB/POZ domain (PFAM PF00651) which is known to occur in TFs [17]. Similarly, (3) the CAD60697

(*Podospora anserine*) is a predicted classical zinc finger protein. (4) The human proteins LZTS1/FEZ1 (Q9Y250, isoform 1) and (5) LZTS2/LAPSER1 (isoform AAK31577) are annotated as leucine zipper putative tumor suppressor 1 and 2, respectively. They harbour a Fez1 domain (PFAM PF06818) which contains a leucine-zipper region with similarity to the DNA-binding domain of the cAMP-responsive activating-transcription factor 5 [18]. There is evidence that Fez1 regulates mitosis and inhibits cancer cell growth [19]. As the accuracy of the protein sequences in the proteomes, especially of the isoforms, improves in the future, more protein examples might pass the selection criteria.

**Supplementary Table 1: Taxonomic distribution of NFAT5 isoforms and their sequence conservation in the N-terminal region of human NFAT5a**

This part contains Supplementary Table 1 (a list of sequence database entries with orthologues of NFAT5 isoforms) and Supplementary Figure 1.1 (an alignment figure of NFAT5 protein isoform sequence segments).

The Supplementary Table 1 lists the entries in non-redundant protein database or UniProt that carry sequences homologous to NFAT5. Generally, protein isoforms are not well studied and it is not surprising that explicit NFAT5a entries are available only for a few species at the time of writing (July 2011); yet, their number has grown since 2004 when the authors have first done this survey.

**Supplementary Table 1**

| Species | NFAT5: isoform a | NFAT5: other isoforms |
|---|---|---|
| *Homo sapiens* (human) | NP_619728.2 | AAD48441.1 (isoform b) O94916 (isoform c) NP_001106649.1 (isoform d) NP_619727.2 BAA74850.2[#] |
| *Pan troglodytes* (chimpanzee) | | XP_001168930.1 |
| *Macaca mulatta* (rhesus monkey) | | XP_001093880.2 |
| *Callithrix jacchus* (white-tufted-ear marmoset) | | XP_002807800.1 |
| *Equus caballus* (horse) | | XP_001497345.2 |
| *Ailuropoda melanoleuca* (giant panda) | | XP_002923698.1 |
| *Canis familiaris* (dog) | | XP_546854.2 |
| *Bos taurus* (cattle) | | XP_002694885.1 |
| *Oryctolagus cuniculus* (rabbit) | | XP_002711712.1 |
| *Mus musculus* (house mouse) | NP_598718.2 | NP_061293.2 AAF31405.1 Q9WV30.1 |
| *Rattus norvegicus* (Norway rat) | EDL92475.1 | NP_001100895.1 EDL92477.1 |
| *Monodelphis domestica* (gray short-tailed opossum) | | XP_001378219.1 |
| *Ornithorhynchus anatinus* (platypus) | | XP_001510046.1 |
| *Gallus gallus* (red jungle fowl) | BAG70407.2 | XP_414226.2 |
| *Xenopus (Silurana) tropicalis* (western clawed frog) | | XP_002937623.1 |

[#] This sequence is annotated as KIAA0827 protein having a length of 1608 AAs. It is an N-terminally prolonged isoform c-type sequence that does not start with methionine. The 77 additional AAs are rich in P (27.9%), R (19.7%), and S (19.7%). Possibly, this is just a hypothetical translation with no real protein equivalent.

## Supplementary Figure 1.1

The alignment figure below (created with CLUSTALX; [20-22]) shows that the protein sequence region in the environment of the human NFAT5a N-terminus is strongly conserved over a wide range of species. The species is encoded with a two-letter abbreviation in accordance with Supplementary Table 1. The N-terminal positions Gly2 and Cys5 are remarkably preserved throughout species.

## CLUSTAL ClustalW 2.0 MULTIPLE SEQUENCE ALIGNMENT

## Supplementary Figure 1.2: Sequences of all five human NFAT5a mRNA transcripts in the regions used for qPCR

The alignment figure below (created with CLUSTALX; [20-22]) shows the mRNA sequences of all five human transcripts in the region of the segments 2|B, -|X and 4|D used for quantitative RT-PCR.

# CLUSTAL ClustalW 2.0 MULTIPLE SEQUENCE ALIGNMENT

```
NM_138713.2|NFAT5d2_t2    ATGCCCTGGACTTCATGTCATTGTGTCAGTGTGGACCTAGACCTGGAATGCCCAAGTCCCTGTACTGTGAGATTGTCTGAAGTTACACCCATCACAGAATTTTGATAGAGTGGACTATTGGAAGAATCTGTGTATGATTTTGTCCCA    150
NM_001113178.1|NFAT5d1_t6 ATGCCCTGGACTTCATGTCATTGCTGAGCGTGGACCTAGACCTGGAATGCCCAAGTCCCTGTACTGGGAGATTGTCTGAAGTTACACCCATCACAGAATTTTGATAGAGTGGACTATTGGAAGAATCTGTGTATGATTTTGTCCCA    150
NM_138714.2|NFAT5a_t1     ATGCCCTGGACTTCATGTCATTGCTGAGCGGGGACCTAGACCTGGAATCGCCCAAGTCCCTGTACTCGCGAGATTGTCTGAAGTTACACCCATCACAGAATTTTGATAGAGTGGACTATTGGAAGAATCTGTGTATGATTTTGTCCCA    150
NM_006599.2|NFAT5c_t3     ATGCCCTGGACTTCATGTCATTGCTGAGCGGGACCTAGACCTGGAATGCCCAAGTCCCTGTACTGGCGAG-----------------------------------------------------AATCTGTGTATGATTTTGTCCCA    96
NM_173214.1|NFAT5a_t4     ATGCCCTGGACTTCATGTCATTGCTGAGCGGGACCTAGACCTGGAATGCCCAAGTCCCTGTACTGGCGAG-----------------------------------------------------AATCTGTGTATGATTTTGTCCCA    96
exon                      AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBCCCCCCCCCCCCCCCCCCCCCCC    150
isoform_a                 -------------------------------------------------------------------------------------------------------------------------------------------------    1
isoform_b                 -------------------------------------------------------------------------------------------------------------------------------------------------    1
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC-----------------------------------------------------CCCCCCCCCCCCCCCCCCCCCCCC    96
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    150
                          1......10........20........30........40........50........60........70........80........90.......100.......110.......120.......130.......140.......150
```

```
NM_138713.2|NFAT5d2_t2    AAGGAGTTACAGTTACCTCCATCTTAGAGAAACATCTGTAGCATCAATGAGTCAGACAAGCGGTGGTGAGGCAGGCTGCCTCCTCCAGCTGTTGTTGTGTGTG----------------------------------------    254
NM_001113178.1|NFAT5d1_t6 AAGGAGTTACAGTTACCTCCATGTTAGAGAAACATCTGTAGCATCAATGAGTCAGACAAGCGGTGGTGAGGCAGGCTGCCTCCTCCAGCTGTTGTTGTGTGTG----------------------------------------    254
NM_138714.2|NFAT5a_t1     AAGGAGTTACAGTTACCTCCATGTTAGAGAAACATCTGTAGCATCAATGAGTCAGACAAGCGGTGGTGAGGCAGGCTGCCTCCTCCAGCTGTTGTTGGTGTGTGGATTTGCCTCTGAAGCAGGGAGTGTCTGCATTAAAAATGACCTGTAG    300
NM_006599.2|NFAT5c_t3     AAGGAGTTACAGTTACCTCCATCTTAGAGAAACATCTGTAGCATCAATGAGTCAGACAAGCGGTGGTGAGGCAGGCTGCCTCCTCCAGCTGTTGTTGTGTGTG----------------------------------------    200
NM_173214.1|NFAT5a_t4     AAGGAGTTACAGTTACCTCCATGTTAGAGAAACATCTGTAGCATCAATGAGTCAGACAAGCGGTGGTGAGGCAGGCTGCCTCCTCCAGCTGTTGTTGGTGTGTGGATTTGCCTCTGAAGCAGGGAGTGTCTGCATTAAAAATGACCTGTAG    246
exon                      CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC----------------------------------------    254
isoform_a                 -------------------------------------------------------------------------------------------------XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX    47
isoform_b                 ------------------------------------------BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB----------------------------------------    59
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC----------------------------------------    200
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD----------------------------------------    254
                          .......160.......170.......180.......190.......200.......210.......220.......230.......240.......250.......260.......270.......280.......290.......300
```

```
NM_138713.2|NFAT5d2_t2    --------------ATGCTTTTTCAGTGTCCCTCCTCTTCCTCCATGGGCGGTGCTTGCAGCTCGTTTACCACCTCTTCCAGCCCTACCATTTATTGTACCTCAGTCACCGACAGCAAGGCTATGCAAGTGGAGAGTTGCTCCTCAGCCG    388
NM_001113178.1|NFAT5d1_t6 --------------ATGCTTTTTCAGTGTCCCTCCTCTTCCTCCATGGGCGGTGCTTGCAGCTCGTTTACCACCTCTTCCAGCCCTACCATTTATTGTACCTCAGTCACCGACAGCAAGGCTATGCAAGTGGAGAGTTGCTCCTCAGCCG    388
NM_138714.2|NFAT5a_t1     TTGTCTGCGTTCATAGATGCTTTTTCAGTGTCCCTCCTCTTCCTCCATGGGCGGTGCTTGCAGCTCGTTTACCACCTCTTCCAGCCCTACCATTTATTGTACCTCAGTCACCGACAGCAAGGCTATGCAAGTGGAGAGTTGCTCCTCAGCCG    450
NM_006599.2|NFAT5c_t3     --------------ATGCTTTTTCAGTGTCCCTCCTCTTCCTCCATGGGCGGTGCTTGCAGCTCGTTTACCACCTCTTCCAGCCCTACCATTTATTGTACCTCAGTCACCGACAGCAAGGCTATGCAAGTGGAGAGTTGCTCCTCAGCCG    334
NM_173214.1|NFAT5a_t4     TTGTCTGCGTTCATAGATGCTTTTTCAGTGTCCCTCCTCTTCCTCCATGGGCGGTGCTTGCAGCTCGTTTACCACCTCTTCCAGCCCTACCATTTATTGTACCTCAGTCACCGACAGCAAGGCTATGCAAGTGGAGAGTTGCTCCTCAGCCG    396
exon                      --------------DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    388
isoform_a                 XXXXXXXXXXXXXX-------------------------AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    168
isoform_b                 --------------BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    193
isoform_c                 --------------CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    334
isoform_d                 --------------DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    388
                          .......310.......320.......330.......340.......350.......360.......370.......380.......390.......400.......410.......420.......430.......440.......450
```

```
NM_138713.2|NFAT5d2_t2    TGGGGGTAAGTAACAGAGGGGTAAGTGAAAAGCAGTTAACCAGTAACACAGTTCAGCAGCATCCATCAACACCGAAGAGGCACACAGTCTTGTACATCTCACCACCACCTGAGGACTTGCTGGATAACAGTTGGATGTCCTGCCAGGATG    538
NM_001113178.1|NFAT5d1_t6 TGGGGGTAAGTAACAGAGGGGTAAGTGAAAAGCAGTTAACCAGTAACACAGTTCAGCAGCATCCATCAACACCGAAGAGGCACACAGTCTTGTACATCTCACCACCACCTGAGGACTTGCTGGATAACAGTTGGATGTCCTGCCAGGATG    538
NM_138714.2|NFAT5a_t1     TGGGGGTAAGTAACAGAGGGGTAAGTGAAAAGCAGTTAACCAGTAACACAGTTCAGCAGCATCCATCAACACCGAAGAGGCACACAGTCTTGTACATCTCACCACCACCTGAGGACTTGCTGGATAACAGTTGGATGTCCTGCCAGGATG    600
NM_006599.2|NFAT5c_t3     TGGGGGTAAGTAACAGAGGGGTAAGTGAAAAGCAGTTAACCAGTAACACAGTTCAGCAGCATCCATCAACACCGAAGAGGCACACAGTCTTGTACATCTCACCACCACCTGAGGACTTGCTGGATAACAGTTGGATGTCCTGCCAGGATG    484
NM_173214.1|NFAT5a_t4     TGGGGGTAAGTAACAGAGGGGTAAGTGAAAAGCAGTTAACCAGTAACACAGTTCAGCAGCATCCATCAACACCGAAGAGGCACACAGTCTTGTACATCTCACCACCACCTGAGGACTTGCTGGATAACAGTTGGATGTCCTGCCAGGATG    546
exon                      DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    538
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    318
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    343
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    484
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    538
                          .......460.......470.......480.......490.......500.......510.......520.......530.......540.......550.......560.......570.......580.......590.......600
```

```
NM_138713.2|NFAT5d2_t2   AGGGGTGTGGATTGGAATCTGAGCAGAGCTGCAGTATGTGGATGGAGGATTCCCCCTCCAACTTCAGTAACATGAGCACCAGTTCCTACAATGATAACACTGAGGTACCTGGTAAATCACGAAAACGAAATCCAAAGCAGAGGCCGGGGG   688
NM_001113178.1|NFAT5d1_t6  AGGGGTGTGGATTGGAATCTGAGCAGAGCTGCAGTATGTGGATGGAGGATTCCCCCTCCAACTTCAGTAACATGAGCACCAGTTCCTACAATGATAACACTGAGGTACCTGGTAAATCACGAAAACGAAATCCAAAGCAGAGGCCGGGGG   688
NM_138714.2|NFAT5a_t1    AGGGGTGTGGATTGGAATCTGAGCAGAGCTGCAGTATGTGGATGGAGGATTCCCCCTCCAACTTCAGTAACATGAGCACCAGTTCCTACAATGATAACACTGAGGTACCTGGTAAATCACGAAAACGAAATCCAAAGCAGAGGCCGGGGG   750
NM_006599.2|NFAT5c_t3    AGGGGTGTGGATTGGAATCTGAGCAGAGCTGCAGTATGTGGATGGAGGATTCCCCCTCCAACTTCAGTAACATGAGCACCAGTTCCTACAATGATAACACTGAGGTACCTGGTAAATCACGAAAACGAAATCCAAAGCAGAGGCCGGGGG   634
NM_173214.1|NFAT5a_t4    AGGGGTGTGGATTGGAATCTGAGCAGAGCTGCAGTATGTGGATGGAGGATTCCCCCTCCAACTTCAGTAACATGAGCACCAGTTCCTACAATGATAACACTGAGGTACCTGGTAAATCACGAAAACGAAATCCAAAGCAGAGGCCGGGGG   696
exon                     DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   688
isoform_a                AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA   468
isoform_b                BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB   493
isoform_c                CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC   634
isoform_d                DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   688
                         .......610.......620.......630.......640.......650.......660.......670.......680.......690.......700.......710.......720.......730.......740.......750
```

```
NM_138713.2|NFAT5d2_t2   TCAAACGACGAGATTGTGAAGAATCTAATATGGATATATTTGATGCCGACAGTGCCAAAGCACCTCACTATGTGCTTTCTCAGCTTACCACGGACAACAAAGGCAACTCAAAAGCGGGAAATGGAACGATTGGAAAACCAAAAAGGAACTG   838
NM_001113178.1|NFAT5d1_t6  TCAAACGACGAGATTGTGAAGAATCTAATATGGATATATTTGATGCCGACAGTGCCAAAGCACCTCACTATGTGCTTTCTCAGCTTACCACGGACAACAAAGGCAACTCAAAAGCGGGAAATGGAACGATTGGAAAACCAAAAAGGAACTG   838
NM_138714.2|NFAT5a_t1    TCAAACGACGAGATTGTGAAGAATCTAATATGGATATATTTGATGCCGACAGTGCCAAAGCACCTCACTATGTGCTTTCTCAGCTTACCACGGACAACAAAGGCAACTCAAAAGCGGGAAATGGAACGATTGGAAAACCAAAAAGGAACTG   900
NM_006599.2|NFAT5c_t3    TCAAACGACGAGATTGTGAAGAATCTAATATGGATATATTTGATGCCGACAGTGCCAAAGCACCTCACTATGTGCTTTCTCAGCTTACCACGGACAACAAAGGCAACTCAAAAGCGGGAAATGGAACGATTGGAAAACCAAAAAGGAACTG   784
NM_173214.1|NFAT5a_t4    TCAAACGACGAGATTGTGAAGAATCTAATATGGATATATTTGATGCCGACAGTGCCAAAGCACCTCACTATGTGCTTTCTCAGCTTACCACGGACAACAAAGGCAACTCAAAAGCGGGAAATGGAACGATTGGAAAACCAAAAAGGAACTG   846
exon                     DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDEEEEEEEEEEEEEEEEEEEE   838
isoform_a                AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA   618
isoform_b                BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB   643
isoform_c                CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC   784
isoform_d                DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   838
                         .......760.......770.......780.......790.......800.......810.......820.......830.......840.......850.......860.......870.......880.......890.......900
```

```
NM_138713.2|NFAT5d2_t2   GAGTAAAGAAGAGCCCTATGTTGTGTGGACAATATCCTGTTAAAAGTGAGGGAAAGGAGCTGAAGATAGTTGTACAACCTGAGACACAGCACCGAGCTCGGTACCTGACTGAGGGCAGCCGTGGCTCAGTGAAAGATAGAACACAGCAAG   988
NM_001113178.1|NFAT5d1_t6  GAGTAAAGAAGAGCCCTATGTTGTGTGGACAATATCCTGTTAAAAGTGAGGGAAAGGAGCTGAAGATAGTTGTACAACCTGAGACACAGCACCGAGCTCGGTACCTGACTGAGGGCAGCCGTGGCTCAGTGAAAGATAGAACACAGCAAG   988
NM_138714.2|NFAT5a_t1    GAGTAAAGAAGAGCCCTATGTTGTGTGGACAATATCCTGTTAAAAGTGAGGGAAAGGAGCTGAAGATAGTTGTACAACCTGAGACACAGCACCGAGCTCGGTACCTGACTGAGGGCAGCCGTGGCTCAGTGAAAGATAGAACACAGCAAG   1050
NM_006599.2|NFAT5c_t3    GAGTAAAGAAGAGCCCTATGTTGTGTGGACAATATCCTGTTAAAAGTGAGGGAAAGGAGCTGAAGATAGTTGTACAACCTGAGACACAGCACCGAGCTCGGTACCTGACTGAGGGCAGCCGTGGCTCAGTGAAAGATAGAACACAGCAAG   934
NM_173214.1|NFAT5a_t4    GAGTAAAGAAGAGCCCTATGTTGTGTGGACAATATCCTGTTAAAAGTGAGGGAAAGGAGCTGAAGATAGTTGTACAACCTGAGACACAGCACCGAGCTCGGTACCTGACTGAGGGCAGCCGTGGCTCAGTGAAAGATAGAACACAGCAAG   996
exon                     EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE   988
isoform_a                AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA   768
isoform_b                BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB   793
isoform_c                CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC   934
isoform_d                DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   988
                         .......910.......920.......930.......940.......950.......960.......970.......980.......990......1000......1010......1020......1030......1040......1050
```

```
NM_138713.2|NFAT5d2_t2   GCTTTCCTACAGTAAAGCTGGAAGGCCATAATGAACCTGTAGTGTTGCAAGTGTTTGTGGGCAACGACTCTGGACGAGTGAAACCACATGGATTTTATCAGGCCTGCAGAGTAACTGGACGAAATACAACTCCTTGCAAAGAAGTGGACA   1138
NM_001113178.1|NFAT5d1_t6  GCTTTCCTACAGTAAAGCTGGAAGGCCATAATGAACCTGTAGTGTTGCAAGTGTTTGTGGGCAACGACTCTGGACGAGTGAAACCACATGGATTTTATCAGGCCTGCAGAGTAACTGGACGAAATACAACTCCTTGCAAAGAAGTGGACA   1138
NM_138714.2|NFAT5a_t1    GCTTTCCTACAGTAAAGCTGGAAGGCCATAATGAACCTGTAGTGTTGCAAGTGTTTGTGGGCAACGACTCTGGACGAGTGAAACCACATGGATTTTATCAGGCCTGCAGAGTAACTGGACGAAATACAACTCCTTGCAAAGAAGTGGACA   1200
NM_006599.2|NFAT5c_t3    GCTTTCCTACAGTAAAGCTGGAAGGCCATAATGAACCTGTAGTGTTGCAAGTGTTTGTGGGCAACGACTCTGGACGAGTGAAACCACATGGATTTTATCAGGCCTGCAGAGTAACTGGACGAAATACAACTCCTTGCAAAGAAGTGGACA   1084
NM_173214.1|NFAT5a_t4    GCTTTCCTACAGTAAAGCTGGAAGGCCATAATGAACCTGTAGTGTTGCAAGTGTTTGTGGGCAACGACTCTGGACGAGTGAAACCACATGGATTTTATCAGGCCTGCAGAGTAACTGGACGAAATACAACTCCTTGCAAAGAAGTGGACA   1146
exon                     EEEEEEEEEEEEEEEEEEEEFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF   1138
isoform_a                AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA   918
isoform_b                BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB   943
isoform_c                CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC   1084
isoform_d                DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   1138
                         ......1060......1070......1080......1090......1100......1110......1120......1130......1140......1150......1160......1170......1180......1190......1200
```

# CLUSTAL ClustalW 2.0 MULTIPLE SEQUENCE ALIGNMENT

```
NM_138713.2|NFAT5d2_t2    CAGCAAGACCTTGCTCTTTTGAAGAGGCCATGAAAGCAATGAAAACTACTGGATGTAATTTAGATAAGGTAAATATTATCCCTAATGCCCTGATGACTCCACTCATACCAAGCAGTATGATTAAGAGTGAAGATGTTACTCCAATGGAAG    1888
NM_001113178.1|NFAT5d1_t6 CAGCAAGACCTTGCTCTTTTGAAGAGGCCATGAAAGCAATGAAAACTACTGGATGTAATTTAGATAAGGTAAATATTATCCCTAATGCCCTGATGACTCCACTCATACCAAGCAGTATGATTAAGAGTGAAGATGTTACTCCAATGGAAG    1885
NM_138714.2|NFAT5a_t1     CAGCAAGACCTTGCTCTTTTGAAGAGGCCATGAAAGCAATGAAAACTACTGGATGTAATTTAGATAAGGTAAATATTATCCCTAATGCCCTGATGACTCCACTCATACCAAGCAGTATGATTAAGAGTGAAGATGTTACTCCAATGGAAG    1950
NM_006599.2|NFAT5c_t3     CAGCAAGACCTTGCTCTTTTGAAGAGGCCATGAAAGCAATGAAAACTACTGGATGTAATTTAGATAAGGTAAATATTATCCCTAATGCCCTGATGACTCCACTCATACCAAGCAGTATGATTAAGAGTGAAGATGTTACTCCAATGGAAG    1834
NM_173214.1|NFAT5a_t4     CAGCAAGACCTTGCTCTTTTGAAGAGGCCATGAAAGCAATGAAAACTACTGGATGTAATTTAGATAAGGTAAATATTATCCCTAATGCCCTGATGACTCCACTCATACCAAGCAGTATGATTAAGAGTGAAGATGTTACTCCAATGGAAG    1896
exon                      LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM    1888
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    1668
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    1693
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    1834
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    1888
                          ......1810......1820......1830......1840......1850......1860......1870......1880......1890......1900......1910......1920......1930......1940......1950
```

```
NM_138713.2|NFAT5d2_t2    TAACAGCAGAAAAAGATCTTCCACTATTTTTAAGACTACAAAGTCTGTTGGATCAACTCAGCAAACATTAGAAAACATCTCAAACATAGCAGGAAATGGCTGTTTTTTCATCACCATCATCTTCCCACCTACCTTCTGAAAATGAAAAC    2038
NM_001113178.1|NFAT5d1_t6 TAACAGCAGAAAAAGATCTTCCACTATTTTTAAGACTACAAAGTCTGTTGGATCAACTCAGCAAACATTAGAAAACATCTCAAACATAGCAGGAAATGGCTGTTTTTTCATCACCATCATCTTCCCACCTACCTTCTGAAAATGAAAAAC    2035
NM_138714.2|NFAT5a_t1     TAACAGCAGAAAAAGATCTTCCACTATTTTTAAGACTACAAAGTCTGTTGGATCAACTCAGCAAACATTAGAAAACATCTCAAACATAGCAGGAAATGGCTGTTTTTTCATCACCATCATCTTCCCACCTACCTTCTGAAAATGAAAAAC    2100
NM_006599.2|NFAT5c_t3     TAACAGCAGAAAAAGATCTTCCACTATTTTTAAGACTACAAAGTCTGTTGGATCAACTCAGCAAACATTAGAAAACATCTCAAACATAGCAGGAAATGGCTGTTTTTTCATCACCATCATCTTCCCACCTACCTTCTGAAATGAAAAC    1984
NM_173214.1|NFAT5a_t4     TAACAGCAGAAAAAGATCTTCCACTATTTTTAAGACTACAAAGTCTGTTGGATCAACTCAGCAAACATTAGAAAACATCTCAAACATAGCAGGAAATGGCTGTTTTTTCATCACCATCATCTTCCCACCTACCTTCTGAAATGAAAAAC    2046
exon                      MMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN    2038
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    1818
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    1843
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    1984
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    2038
                          ......1960......1970......1980......1990......2000......2010......2020......2030......2040......2050......2060......2070......2080......2090......2100
```

```
NM_138713.2|NFAT5d2_t2    AGCAGCAGATTCAGCCCAAGGCATACAACCCAGAGACCCTGACAACTATTCAAACCCAGGACATCTCACAGCCTGGTACTTTTTCCAGCAGTTTGTCTTCTAGTCAGCTGCCCAAGACGCATGCACTATTGCAGCAGGCTACACAGTTTC    2188
NM_001113178.1|NFAT5d1_t6 AGCAGCAGATTCAGCCCAAGGCATACAACCCAGAGACCCTGACAACTATTCAAACCCAGGACATCTCACAGCCTGGTACTTTTTCCAGCAGTTTGTCTTCTAGTCAGCTGCCCAAGACGCATGCACTATTGCAGCAGGCTACACAGTTTC    2185
NM_138714.2|NFAT5a_t1     AGCAGCAGATTCAGCCCAAGGCATACAACCCAGAGACCCTGACAACTATTCAAACCCAGGACATCTCACAGCCTGGTACTTTTTCCAGCAGTTTGTCTTCTAGTCAGCTGCCCAAGACGCATGCACTATTGCAGCAGGCTACACAGTTTC    2250
NM_006599.2|NFAT5c_t3     AGCAGCAGATTCAGCCCAAGGCATACAACCCAGAGACCCTGACAACTATTCAAACCCAGGACATCTCACAGCCTGGTACTTTTTCCAGCAGTTTGTCTTCTAGTCAGCTGCCCAAGACGCATGCACTATTGCAGCAGGCTACACAGTTTC    2134
NM_173214.1|NFAT5a_t4     AGCAGCAGATTCAGCCCAAGGCATACAACCCAGAGACCCTGACAACTATTCAAACCCAGGACATCTCACAGCCTGGTACTTTTTCCAGCAGTTTGTCTTCTAGTCAGCTGCCCAAGACGCATGCACTATTGCAGCAGGCTACACAGTTTC    2196
exon                      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN    2188
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    1968
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    1993
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    2134
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    2188
                          ......2110......2120......2130......2140......2150......2160......2170......2180......2190......2200......2210......2220......2230......2240......2250
```

```
NM_138713.2|NFAT5d2_t2    AGACAAGAGAAACTCAGTCTAGAGAGATATTACAGTCAGATGGTACAGTGGTTAATTTGTCACAACTGACTGAGGCATCACAACAACAGCAGCAGTCACCACTACAAGAACAAGCACAGACTTTACAGCAGCAGATTTCATCAAATATTT    2338
NM_001113178.1|NFAT5d1_t6 AGACAAGAGAAACTCAGTCTAGAGAGATATTACAGTCAGATGGTACAGTGGTTAATTTGTCACAACTGACTGAGGCATCACAACAACAGCAGCAGTCACCACTACAAGAACAAGCACAGACTTTACAGCAGCAGATTTCATCAAATATTT    2335
NM_138714.2|NFAT5a_t1     AGACAAGAGAAACTCAGTCTAGAGAGATATTACAGTCAGATGGTACAGTGGTTAATTTGTCACAACTGACTGAGGCATCACAACAACAGCAGCAGTCACCACTACAAGAACAAGCACAGACTTTACAGCAGCAGATTTCATCAAATATTT    2400
NM_006599.2|NFAT5c_t3     AGACAAGAGAAACTCAGTCTAGAGAGATATTACAGTCAGATGGTACAGTGGTTAATTTGTCACAACTGACTGAGGCATCACAACAACAGCAGCAGTCACCACTACAAGAACAAGCACAGACTTTACAGCAGCAGATTTCATCAAATATTT    2284
NM_173214.1|NFAT5a_t4     AGACAAGAGAAACTCAGTCTAGAGAGATATTACAGTCAGATGGTACAGTGGTTAATTTGTCACAACTGACTGAGGCATCACAACAACAGCAGCAGTCACCACTACAAGAACAAGCACAGACTTTACAGCAGCAGATTTCATCAAATATTT    2346
exon                      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN    2338
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    2118
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    2143
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    2284
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    2338
                          ......2260......2270......2280......2290......2300......2310......2320......2330......2340......2350......2360......2370......2380......2390......2400
```

```
NM_138713.2|NFAT5d2_t2    TTCCATCACCAAATAGTGTGAGTCAGCTTCAGAATACTATTCAGCAGCTGCAAGCAGGGAGTTTCACAGGCAGTACTGCTAGTGGCAGCAGTGGAAGTGTTGACTTGGTCCAACAAGTTTTAGAGGCACAGCAGCAGTTATCTTCAGTTT    2488
NM_001113178.1|NFAT5d1_t6 TTCCATCACCAAATAGTGTGAGTCAGCTTCAGAATACTATTCAGCAGCTGCAAGCAGGGAGTTTCACAGGCAGTACTGCTAGTGGCAGCAGTGGAAGTGTTGACTTGGTCCAACAAGTTTTAGAGGCACAGCAGCAGTTATCTTCAGTTT    2485
NM_138714.2|NFAT5a_t1     TTCCATCACCAAATAGTGTGAGTCAGCTTCAGAATACTATTCAGCAGCTGCAAGCAGGGAGTTTCACAGGCAGTACTGCTAGTGGCAGCAGTGGAAGTGTTGACTTGGTCCAACAAGTTTTAGAGGCACAGCAGCAGTTATCTTCAGTTT    2550
NM_006599.2|NFAT5c_t3     TTCCATCACCAAATAGTGTGAGTCAGCTTCAGAATACTATTCAGCAGCTGCAAGCAGGGAGTTTCACAGGCAGTACTGCTAGTGGCAGCAGTGGAAGTGTTGACTTGGTCCAACAAGTTTTAGAGGCACAGCAGCAGTTATCTTCAGTTT    2434
NM_173214.1|NFAT5a_t4     TTCCATCACCAAATAGTGTGAGTCAGCTTCAGAATACTATTCAGCAGCTGCAAGCAGGGAGTTTCACAGGCAGTACTGCTAGTGGCAGCAGTGGAAGTGTTGACTTGGTCCAACAAGTTTTAGAGGCACAGCAGCAGTTATCTTCAGTTT    2496
exon                      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN    2488
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    2268
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    2293
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    2434
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    2488
                          ......2410......2420......2430......2440......2450......2460......2470......2480......2490......2500......2510......2520......2530......2540......2550
```

```
NM_138713.2|NFAT5d2_t2    TATTTTTGCTCCAGATGGTAATGAGAATGTTCAAGAGCAGCTTAGTGCAGATATTTTTCAACAAGTCAGTCAAATTCAGAGTGGTGTAAGCCCTGGAATGTTTTCCTCAACAGAGCCAACAGTCCATACCGACCAGATAATTTATTAC    2638
NM_001113178.1|NFAT5d1_t6 TATTTTTGCTCCAGATGGTAATGAGAATGTTCAAGAGCAGCTTAGTGCAGATATTTTTCAACAAGTCAGTCAAATTCAGAGTGGTGTAAGCCCTGGAATGTTTTCCTCAACAGAGCCAACAGTCCATACCGACCAGATAATTTATTAC    2635
NM_138714.2|NFAT5a_t1     TATTTTTGCTCCAGATGGTAATGAGAATGTTCAAGAGCAGCTTAGTGCAGATATTTTTCAACAAGTCAGTCAAATTCAGAGTGGTGTAAGCCCTGGAATGTTTTCCTCAACAGAGCCAACAGTCCATACCGACCAGATAATTTATTAC    2700
NM_006599.2|NFAT5c_t3     TATTTTTGCTCCAGATGGTAATGAGAATGTTCAAGAGCAGCTTAGTGCAGATATTTTTCAACAAGTCAGTCAAATTCAGAGTGGTGTAAGCCCTGGAATGTTTTCCTCAACAGAGCCAACAGTCCATACCGACCAGATAATTTATTAC    2584
NM_173214.1|NFAT5a_t4     TATTTTTGCTCCAGATGGTAATGAGAATGTTCAAGAGCAGCTTAGTGCAGATATTTTTCAACAAGTCAGTCAAATTCAGAGTGGTGTAAGCCCTGGAATGTTTTCCTCAACAGAGCCAACAGTCCATACCGACCAGATAATTTATTAC    2646
exon                      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN    2638
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    2418
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    2443
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    2584
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    2638
                          ......2560......2570......2580......2590......2600......2610......2620......2630......2640......2650......2660......2670......2680......2690......2700
```

```
NM_138713.2|NFAT5d2_t2    CTGGAAGAGCTGAAAGTGTTCATCCACAGTCTGAAAACAGCGTTATCTAATAACAGCAGCAGCAGCAGCAACAGCAAGTGATGGAATCTTCAGCCGCCAATGGTGATGGAGATGCAACAGAGTATCTGCTCAGGCAGGTGCCCAGATTC    2788
NM_001113178.1|NFAT5d1_t6 CTGGAAGAGCTGAAAGTGTTCATCCACAGTCTGAAAACAGCGTTATCTAATAACAGCAGCAGCAGCAGCAACAGCAAGTGATGGAATCTTCAGCCGCCAATGGTGATGGAGATGCAACAGAGTATCTGCTCAGGCAGGTGCCCAGATTC    2785
NM_138714.2|NFAT5a_t1     CTGGAAGAGCTGAAAGTGTTCATCCACAGTCTGAAAACAGCGTTATCTAATAACAGCAGCAGCAGCAGCAACAGCAAGTGATGGAATCTTCAGCCGCCAATGGTGATGGAGATGCAACAGAGTATCTGCTCAGGCAGGTGCCCAGATTC    2850
NM_006599.2|NFAT5c_t3     CTGGAAGAGCTGAAAGTGTTCATCCACAGTCTGAAAACAGCGTTATCTAATAACAGCAGCAGCAGCAGCAACAGCAAGTGATGGAATCTTCAGCCGCCAATGGTGATGGAGATGCAACAGAGTATCTGCTCAGGCAGGTGCCCAGATTC    2734
NM_173214.1|NFAT5a_t4     CTGGAAGAGCTGAAAGTGTTCATCCACAGTCTGAAAACAGCGTTATCTAATAACAGCAGCAGCAGCAGCAACAGCAAGTGATGGAATCTTCAGCCGCCAATGGTGATGGAGATGCAACAGAGTATCTGCTCAGGCAGGTGCCCAGATTC    2796
exon                      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN    2788
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    2568
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    2593
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    2734
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    2788
                          ......2710......2720......2730......2740......2750......2760......2770......2780......2790......2800......2810......2820......2830......2840......2850
```

```
NM_138713.2|NFAT5d2_t2    AGTCAGAGTTATTCCCTTCAACTGCTTCAGCAAATGGAAACCTTCAGCAATGCCAGTTTACCAGCAGACTTCTCACATGATGAGTGCATTGTCTACCAATGAGGATATGCAAATGCAGTGTGAATTGTTTTCTTCTCCTCCTGCAGTTT    2938
NM_001113178.1|NFAT5d1_t6 AGTCAGAGTTATTCCCTTCAACTGCTTCAGCAAATGGAAACCTTCAGCAATGCCAGTTTACCAGCAGACTTCTCACATGATGAGTGCATTGTCTACCAATGAGGATATGCAAATGCAGTGTGAATTGTTTTCTTCTCCTCCTGCAGTTT    2935
NM_138714.2|NFAT5a_t1     AGTCAGAGTTATTCCCTTCAACTGCTTCAGCAAATGGAAACCTTCAGCAATGCCAGTTTACCAGCAGACTTCTCACATGATGAGTGCATTGTCTACCAATGAGGATATGCAAATGCAGTGTGAATTGTTTTCTTCTCCTCCTGCAGTTT    3000
NM_006599.2|NFAT5c_t3     AGTCAGAGTTATTCCCTTCAACTGCTTCAGCAAATGGAAACCTTCAGCAATGCCAGTTTACCAGCAGACTTCTCACATGATGAGTGCATTGTCTACCAATGAGGATATGCAAATGCAGTGTGAATTGTTTTCTTCTCCTCCTGCAGTTT    2884
NM_173214.1|NFAT5a_t4     AGTCAGAGTTATTCCCTTCAACTGCTTCAGCAAATGGAAACCTTCAGCAATGCCAGTTTACCAGCAGACTTCTCACATGATGAGTGCATTGTCTACCAATGAGGATATGCAAATGCAGTGTGAATTGTTTTCTTCTCCTCCTGCAGTTT    2946
exon                      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN    2938
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    2718
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    2743
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    2884
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    2938
                          ......2860......2870......2880......2890......2900......2910......2920......2930......2940......2950......2960......2970......2980......2990......3000
```

```
NM_138713.2|NFAT5d2_t2   CTGGAAATGAAACTTCTACAACTACCACACAGCAGGTTGCAACCCCTGGCACTACCATGTTTCAGACATCAAGTTCAGGAGATGGAGAAGAAACTGGAACACAAGCAAAACAGATTCAGAACAGTGTCTTTCAGACCATGGTCCAAATGC   3088
NM_001113178.1|NFAT5d1_t6 CTGGAAATGAAACTTCTACAACTACCACACAGCAGGTTGCAACCCCTGGCACTACCATGTTTCAGACATCAAGTTCAGGAGATGGAGAAGAAACTGGAACACAAGCAAAACAGATTCAGAACAGTGTCTTTCAGACCATGGTCCAAATGC   3085
NM_138714.2|NFAT5a_t1    CTGGAAATGAAACTTCTACAACTACCACACAGCAGGTTGCAACCCCTGGCACTACCATGTTTCAGACATCAAGTTCAGGAGATGGAGAAGAAACTGGAACACAAGCAAAACAGATTCAGAACAGTGTCTTTCAGACCATGGTCCAAATGC   3150
NM_006599.2|NFAT5c_t3    CTGGAAATGAAACTTCTACAACTACCACACAGCAGGTTGCAACCCCTGGCACTACCATGTTTCAGACATCAAGTTCAGGAGATGGAGAAGAAACTGGAACACAAGCAAAACAGATTCAGAACAGTGTCTTTCAGACCATGGTCCAAATGC   3034
NM_173214.1|NFAT5a_t4    CTGGAAATGAAACTTCTACAACTACCACACAGCAGGTTGCAACCCCTGGCACTACCATGTTTCAGACATCAAGTTCAGGAGATGGAGAAGAAACTGGAACACAAGCAAAACAGATTCAGAACAGTGTCTTTCAGACCATGGTCCAAATGC   3096
exon                     NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN   3088
isoform_a                AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA   2868
isoform_b                BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB   2893
isoform_c                CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC   3034
isoform_d                DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   3088
                         ......3010......3020......3030......3040......3050......3060......3070......3080......3090......3100......3110......3120......3130......3140......3150
```

```
NM_138713.2|NFAT5d2_t2   AACATAGTGGGGACAATCAACCTCAAGTTAACCTTTTTTCATCCACAAAAAGTATGATGAGTGTTCAGAATAGTGGTACCCAACAACAAGGTAATGGTTTATTCCAGCAAGGGAATGAGATGATGTCACTTCAATCTGGAAATTTTTTGC   3238
NM_001113178.1|NFAT5d1_t6 AACATAGTGGGGACAATCAACCTCAAGTTAACCTTTTTTCATCCACAAAAAGTATGATGAGTGTTCAGAATAGTGGTACCCAACAACAAGGTAATGGTTTATTCCAGCAAGGGAATGAGATGATGTCACTTCAATCTGGAAATTTTTTGC   3235
NM_138714.2|NFAT5a_t1    AACATAGTGGGGACAATCAACCTCAAGTTAACCTTTTTTCATCCACAAAAAGTATGATGAGTGTTCAGAATAGTGGTACCCAACAACAAGGTAATGGTTTATTCCAGCAAGGGAATGAGATGATGTCACTTCAATCTGGAAATTTTTTGC   3300
NM_006599.2|NFAT5c_t3    AACATAGTGGGGACAATCAACCTCAAGTTAACCTTTTTTCATCCACAAAAAGTATGATGAGTGTTCAGAATAGTGGTACCCAACAACAAGGTAATGGTTTATTCCAGCAAGGGAATGAGATGATGTCACTTCAATCTGGAAATTTTTTGC   3184
NM_173214.1|NFAT5a_t4    AACATAGTGGGGACAATCAACCTCAAGTTAACCTTTTTTCATCCACAAAAAGTATGATGAGTGTTCAGAATAGTGGTACCCAACAACAAGGTAATGGTTTATTCCAGCAAGGGAATGAGATGATGTCACTTCAATCTGGAAATTTTTTGC   3246
exon                     NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN   3238
isoform_a                AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA   3018
isoform_b                BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB   3043
isoform_c                CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC   3184
isoform_d                DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   3238
                         ......3160......3170......3180......3190......3200......3210......3220......3230......3240......3250......3260......3270......3280......3290......3300
```

```
NM_138713.2|NFAT5d2_t2   AGCAGTGTTTCTCATTCACAGGCCCAACTTTTTCATCCTCAAAATCCTATTGCCGATGCTCAGAACGTTTCCCAGGAAATCAAGGTTTCTCTTTTCATAGTCCAAATCTATTGTCCACAGTCGAGATTCTACAACCTCCTTTGAACAAA   3388
NM_001113178.1|NFAT5d1_t6 AGCAGTGTTTCTCATTCACAGGCCCAACTTTTTCATCCTCAAAATCCTATTGCCGATGCTCAGAACGTTTCCCAGGAAACTCAAGGTTGCTCTGTTTCATAGTCCAAATCTATTGTCCACAGTCGAGATTCTACAACCTCCTTTGAACAAA   3385
NM_138714.2|NFAT5a_t1    AGCAGTGTTTCTCATTCACAGGCCCAACTTTTTCATCCTCAAAATCCTATTGCCGATGCTCAGAACGTTTCCCAGGAAATCAAGGTTTCTCTTTTCATAGTCCAAATCTATTGTCCACAGTCGAGATTCTACAACCTCCTTTGAACAAA   3450
NM_006599.2|NFAT5c_t3    AGCAGTGTTTCTCATTCACAGGCCCAACTTTTTCATCCTCAAAATCCTATTGCCGATGCTCAGAACGTTTCCCAGGAAATCAAGGTTTCTCTTTTCATAGTCCAAATCTATTGTCCACAGTCGAGATTCTACAACCTCCTTTGAACAAA   3334
NM_173214.1|NFAT5a_t4    AGCAGTGTTTCTCATTCACAGGCCCAACTTTTTCATCCTCAAAATCCTATTGCCGATGCTCAGAACGTTTCCCAGGAAATCAAGGTTTCTCTTTTCATAGTCCAAATCTATTGTCCACAGTCGAGATTCTACAACCTCCTTTGAACAAA   3396
exon                     NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN   3388
isoform_a                AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA   3168
isoform_b                BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB   3193
isoform_c                CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC   3334
isoform_d                DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   3388
                         ......3310......3320......3330......3340......3350......3360......3370......3380......3390......3400......3410......3420......3430......3440......3450
```

```
NM_138713.2|NFAT5d2_t2   TGCAGCCTCCAATGTTTCACTTTCAAAGTACCATTGTGTGTTACAGGGCTTTCAGTTCCTCAAGACCAGCAGTCAACCAACATATTTCTTTCCCAGAGTCCATGAATAATCTTCAGACTAACACAGTAGCCCAAGAAGCATTTTTTG   3538
NM_001113178.1|NFAT5d1_t6 TGCAGCCTCCAATGTTTCACTTTCAAAGTACCATTGTGTGTTACAGGGCTTTCAGTTCCTCAAGACCAGCAGTCAACCAACATATTTCTTTCCCAGAGTCCATGAATAATCTTCAGACTAACACAGTAGCCCAAGAAGCATTTTTTTG   3535
NM_138714.2|NFAT5a_t1    TGCAGCCTCCAATGTTTCACTTTCAAAGTACCATTGTGTGTTACAGGGCTTTCAGTTCCTCAAGACCAGCAGTCAACCAACATATTTCTTTCCCAGAGTCCATGAATAATCTTCAGACTAACACAGTAGCCCAAGAAGCATTTTTTG   3600
NM_006599.2|NFAT5c_t3    TGCAGCCTCCAATGTTTCACTTTCAAAGTACCATTGTGTGTTACAGGGCTTTCAGTTCCTCAAGACCAGCAGTCAACCAACATATTTCTTTCCCAGAGTCCATGAATAATCTTCAGACTAACACAGTAGCCCAAGAAGCATTTTTTG   3484
NM_173214.1|NFAT5a_t4    TGCAGCCTCCAATGTTTCACTTTCAAAGTACCATTGTGTGTTACAGGGCTTTCAGTTCCTCAAGACCAGCAGTCAACCAACATATTTCTTTCCCAGAGTCCATGAATAATCTTCAGACTAACACAGTAGCCCAAGAAGCATTTTTTG   3546
exon                     NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN   3538
isoform_a                AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA   3318
isoform_b                BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB   3343
isoform_c                CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC   3484
isoform_d                DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD   3538
                         ......3460......3470......3480......3490......3500......3510......3520......3530......3540......3550......3560......3570......3580......3590......3600
```

```
NM_138713.2|NFAT5d2_t2    AAGGGTCACCTAGTTCTCAAGAGCAGCAAGTAACTCTCTTCTTTATCTCCAGCATCCATGTCTGCCTTGCAGACCAGTATAAATCAACAAGATATGCAACAGTCTCCTCTTTATTCCCCTCAGAACAACATGCCTGGAATTCAAGGAGCCA    4288
NM_001113178.1|NFAT5d1_t6 AAGGGTCACCTAGTTCTCAAGAGCAGCAAGTAACTCTCTTCTTTATCTCCAGCATCCATGTCTGCCTTGCAGACCAGTATAAATCAACAAGATATGCAACAGTCTCCTCTTTATTCCCCTCAGAACAACATGCCTGGAATTCAAGGAGCCA    4285
NM_138714.2|NFAT5a_t1     AAGGGTCACCTAGTTCTCAAGAGCAGCAAGTAACTCTCTTCTTTATCTCCAGCATCCATGTCTGCCTTGCAGACCAGTATAAATCAACAAGATATGCAACAGTCTCCTCTTTATTCCCCTCAGAACAACATGCCTGGAATTCAAGGAGCCA    4350
NM_006599.2|NFAT5c_t3     AAGGGTCACCTAGTTCTCAAGAGCAGCAAGTAACTCTCTTCTTTATCTCCAGCATCCATGTCTGCCTTGCAGACCAGTATAAATCAACAAGATATGCAACAGTCTCCTCTTTATTCCCCTCAGAACAACATGCCTGGAATTCAAGGAGCCA    4234
NM_173214.1|NFAT5a_t4     AAGGGTCACCTAGTTCTCAAGAGCAGCAAGTAACTCTCTTCTTTATCTCCAGCATCCATGTCTGCCTTGCAGACCAGTATAAATCAACAAGATATGCAACAGTCTCCTCTTTATTCCCCTCAGAACAACATGCCTGGAATTCAAGGAGCCA    4296
exon                      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN    4288
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    4068
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    4093
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    4234
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    4288
                          ......4210......4220......4230......4240......4250......4260......4270......4280......4290......4300......4310......4320......4330......4340......4350


NM_138713.2|NFAT5d2_t2    CATCTTTGCCTCAACCACAGGTTACTTTATTTCACAACACAGCAGGAGGCACAATGAACCAACTGCAGAATTGTCCTGGCTCATCTCAGCAGACATCAGGAATGTTCTTATTTGGCATTCAAAATAACTGTAGTCAGCTTTTAACCTCTTG    4438
NM_001113178.1|NFAT5d1_t6 CATCTTTGCCTCAACCACAGGTTACTTTATTTCACAACACAGCAGGAGGCACAATGAACCAACTGCAGAATTGTCCTGGCTCATCTCAGCAGACATCAGGAATGTTCTTATTTGGCATTCAAAATAACTGTAGTCAGCTTTTAACCTCTTG    4435
NM_138714.2|NFAT5a_t1     CATCTTTGCCTCAACCACAGGTTACTTTATTTCACAACACAGCAGGAGGCACAATGAACCAACTGCAGAATTGTCCTGGCTCATCTCAGCAGACATCAGGAATGTTCTTATTTGGCATTCAAAATAACTGTAGTCAGCTTTTAACCTCTTG    4500
NM_006599.2|NFAT5c_t3     CATCTTTGCCTCAACCACAGGTTACTTTATTTCACAACACAGCAGGAGGCACAATGAACCAACTGCAGAATTGTCCTGGCTCATCTCAGCAGACATCAGGAATGTTCTTATTTGGCATTCAAAATAACTGTAGTCAGCTTTTAACCTCTTG    4384
NM_173214.1|NFAT5a_t4     CATCTTTGCCTCAACCACAGGTTACTTTATTTCACAACACAGCAGGAGGCACAATGAACCAACTGCAGAATTGTCCTGGCTCATCTCAGCAGACATCAGGAATGTTCTTATTTGGCATTCAAAATAACTGTAGTCAGCTTTTAACCTCTTG    4446
exon                      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNPPPPPPPPPPPPPPPPPPPPPP    4438
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    4218
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    4243
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    4384
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    4438
                          ......4360......4370......4380......4390......4400......4410......4420......4430......4440......4450......4460......4470......4480......4490......4500


NM_138713.2|NFAT5d2_t2    GACCAGTTACATTGCCTGATCAGTTGATGGCCATAAGTCAGCCAGGCCAACCACAAAACGAGGGCCAGCCACCTGTGACAACATTGTTTTCTCAGCAAATGCCAGAGAATTGTCCAGTGGCATCTTTATAAACACCAACAGAACACATCG    4588
NM_001113178.1|NFAT5d1_t6 GACCAGTTACATTGCCTGATCAGTTGATGGCCATAAGTCAGCCAGGCCAACCACAAAACGAGGGCCAGCCACCTGTGACAACACTTGTTTTCTCAGCAAATGCCAGAGAATTGTCCAGTGGCATCCTTTATAAACACCAACAGAACATCG    4585
NM_138714.2|NFAT5a_t1     GACCAGTTACATTGCCTGATCAGTTGATGGCCATAAGTCAGCCAGGCCAACCACAAAACGAGGGCCAGCCACCTGTGACAACACTTGTTTTCTCAGCAAATGCCAGAGAATTGTCCAGTGGCATCCTTTATAAACACCAACAGAACATCG    4650
NM_006599.2|NFAT5c_t3     GACCAGTTACATTGCCTGATCAGTTGATGGCCATAAGTCAGCCAGGCCAACCACAAAACGAGGGCCAGCCACCTGTGACAACACTTGTTTTCTCAGCAAATGCCAGAGAATTGTCCAGTGGCATCCTTTATAAACACCAACAGAACATCG    4534
NM_173214.1|NFAT5a_t4     GACCAGTTACATTGCCTGATCAGTTGATGGCCATAAGTCAGCCAGGCCAACCACAAAACGAGGGCCAGCCACCTGTGACAACACTTGTTTTCTCAGCAAATGCCAGAGAATTGTCCAGTGGCATCCTTTATAAACACCAACAGAACATCG    4596
exon                      PPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPP    4588
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    4368
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    4393
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    4534
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    4588
                          ......4510......4520......4530......4540......4550......4560......4570......4580......4590......4600......4610......4620......4630......4640......4650


NM_138713.2|NFAT5d2_t2    AAAAGATTGATTTGCTTGTTTCATTGCAAAACCAAGGGAACAACTTGACTGGCTCCTTTTAA    4650
NM_001113178.1|NFAT5d1_t6 AAAAGATTGATTTGCTTGTTTCATTGCAAAACCAAGGGAACAACTTGACTGGCTCCTTTTAA    4647
NM_138714.2|NFAT5a_t1     AAAAGATTGATTTGCTTGTTTCATTGCAAAACCAAGGGAACAACTTGACTGGCTCCTTTTAA    4712
NM_006599.2|NFAT5c_t3     AAAAGATTGATTTGCTTGTTTCATTGCAAAACCAAGGGAACAACTTGACTGGCTCCTTTTAA    4596
NM_173214.1|NFAT5a_t4     AAAAGATTGATTTGCTTGTTTCATTGCAAAACCAAGGGAACAACTTGACTGGCTCCTTTTAA    4658
exon                      PPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPP    4650
isoform_a                 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA    4430
isoform_b                 BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB    4455
isoform_c                 CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC    4596
isoform_d                 DDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDDD    4650
                          ......4660......4670......4680......4690......4700......4710..
```

*References*

References

1. Eddy SR. What is a hidden Markov model? *Nat Biotechnol*. 2004;22:1315-1316.

2. Finn RD, Tate J, Mistry J et al. The Pfam protein families database. *Nucleic Acids Res*. 2008;36:D281-D288.

3. Letunic I, Copley RR, Pils B et al. SMART 5: domains in the context of genomes and networks. *Nucleic Acids Res*. 2006;34:D257-D260.

4. Maurer-Stroh S, Eisenhaber B, Eisenhaber F. N-terminal N-myristoylation of proteins: prediction of substrate proteins from amino acid sequence. *J Mol Biol*. 2002;317:541-557.

5. Maurer-Stroh S, Gouda M, Novatchkova M et al. MYRbase: analysis of genome-wide glycine myristoylation enlarges the functional spectrum of eukaryotic myristoylated proteins. *Genome Biol*. 2004;5:R21.

6. Maurer-Stroh S, Eisenhaber F. Refinement and prediction of protein prenylation motifs. *Genome Biol*. 2005;6:R55.

7. Maurer-Stroh S, Koranda M, Benetka W et al. Towards complete sets of farnesylated and geranylgeranylated proteins. *PLoS Comput Biol*. 2007;3:e66.

8. Eisenhaber B, Bork P, Eisenhaber F. Prediction of potential GPI-modification sites in proprotein sequences. *J Mol Biol*. 1999;292:741-758.

9. Lopez-Rodriguez C, Aramburu J, Rakeman AS et al. NFAT5, a constitutively nuclear NFAT protein that does not cooperate with Fos and Jun. *Proc Natl Acad Sci U S A*. 1999;96:7214-7219.

10. Miyakawa H, Woo SK, Dahl SC et al. Tonicity-responsive enhancer binding protein, a rel-like protein that stimulates transcription in response to hypertonicity. *Proc Natl Acad Sci U S A*. 1999;96:2538-2542.

11. Pan S, Tsuruta R, Masuda ES et al. NFATz: a novel rel similarity domain containing protein. *Biochem Biophys Res Commun*. 2000;272:765-776.

12. Trama J, Lu Q, Hawley RG et al. The NFAT-related protein NFATL1 (TonEBP/NFAT5) is induced upon T cell activation in a calcineurin-dependent manner. *J Immunol*. 2000;165:4884-4894.

13. Ko BC, Ruepp B, Bohren KM et al. Identification and characterization of multiple osmotic response sequences in the human aldose reductase gene. *J Biol Chem*. 1997;272:16431-16437.

14. Tong EH, Guo JJ, Huang AL et al. Regulation of nucleocytoplasmic trafficking of transcription factor OREBP/TonEBP/NFAT5. *J Biol Chem*. 2006;281:23870-23879.

15. Hogenesch JB, Gu YZ, Moran SM et al. The basic helix-loop-helix-PAS protein MOP9 is a brain-specific heterodimeric partner of circadian and hypoxia factors. *J Neurosci*. 2000;20:RC83.

16. Navarro-Lerida I, varez-Barrientos A, Gavilanes F et al. Distance-dependent cellular palmitoylation of de-novo-designed sequences and their translocation to plasma membrane subdomains. *J Cell Sci*. 2002;115:3119-3130.

17. Stogios PJ, Downs GS, Jauhal JJ et al. Sequence and structural analysis of BTB domain proteins. *Genome Biol*. 2005;6:R82.

18. Ishii H, Baffa R, Numata SI et al. The FEZ1 gene at chromosome 8p22 encodes a leucine-zipper protein, and its expression is altered in multiple human tumors. *Proc Natl Acad Sci U S A*. 1999;96:3928-3933.

19. Ishii H, Vecchione A, Murakumo Y et al. FEZ1/LZTS1 gene at 8p22 suppresses cancer cell growth and regulates mitosis. *Proc Natl Acad Sci U S A*. 2001;98:10374-10379.

20. Thompson JD, Higgins DG, Gibson TJ. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res*. 1994;22:4673-4680.

21. Thompson JD, Gibson TJ, Higgins DG. Multiple sequence alignment using ClustalW and ClustalX. *Curr Protoc Bioinformatics*. 2002;Chapter 2:Unit.

22. Larkin MA, Blackshields G, Brown NP et al. Clustal W and Clustal X version 2.0. *Bioinformatics*. 2007;23:2947-2948.