

SOFTWARE

Open Access

Scan for Motifs: a webserver for the analysis of post-transcriptional regulatory elements in the 3' untranslated regions (3' UTRs) of mRNAs

Ambarish Biswas¹ and Chris M Brown^{1,2*}

Abstract

Background: Gene expression in vertebrate cells may be controlled post-transcriptionally through regulatory elements in mRNAs. These are usually located in the untranslated regions (UTRs) of mRNA sequences, particularly the 3'UTRs.

Results: Scan for Motifs (SFM) simplifies the process of identifying a wide range of regulatory elements on alignments of vertebrate 3'UTRs. SFM includes identification of both RNA Binding Protein (RBP) sites and targets of miRNAs. In addition to searching pre-computed alignments, the tool provides users the flexibility to search their own sequences or alignments. The regulatory elements may be filtered by expected value cutoffs and are cross-referenced back to their respective sources and literature. The output is an interactive graphical representation, highlighting potential regulatory elements and overlaps between them. The output also provides simple statistics and links to related resources for complementary analyses. The overall process is intuitive and fast. As SFM is a free web-application, the user does not need to install any software or databases.

Conclusions: Visualisation of the binding sites of different classes of effectors that bind to 3'UTRs will facilitate the study of regulatory elements in 3' UTRs.

Keywords: Untranslated region, microRNA, RNA binding protein, Translational control

Background

The untranslated regions of mRNA sequences (UTRs) include most of the experimentally determined regulatory elements (REs) [1,2]. This post-transcriptional regulatory information can affect the site at which a mRNA is polyadenylated, and then how, when and where it is translated [3,4]. A number of tools and methods have been developed to identify *cis*-regulatory elements (CREs), many focusing on individual types of CREs in single sequences [5,6]. These may ignore the detection of other types of CREs in the neighboring regions [7,8]. For example, although there are a large number of algorithms to predict microRNA (miRNA) binding sites, reviewed in [9,10], only one has included specific consideration of a nearby RNA binding protein (RBP) site [11]. However, some miRNA targets are known to be affected by the presence of other elements or sequences nearby

[1,11-13]. Most regulatory elements are quite small (<12 bases) and many *in silico* predictions have high false positive rates. Visualisation of potential sites could improve the utility of predictions.

Some complex RNA elements can be both miRNA target sites and be bound by proteins [3,14,15]. Recent publications have shown evidence that specific types of miRNAs and RBPs work in concert to influence transcript decay [11,16,17] or translation [13] and this synergy has been included in some computational analyses for proteins [18] and miRNAs [19].

In many studies one specific gene of interest from a single species is being analysed. Recently developed systems: RegRNA 2.0 [2], AURA [20], ARESite [6], and UTRdb [21] have provided increasing support for this type of analysis. However, the analysis of sequence alignments, a representation of overlapping identified elements, E-value cutoff, and the ability to include custom sequence motifs in the analysis, are not currently available in a single tool. Scan for Motifs provides this for 3'UTR regions. It is primarily aimed at the analysis of

* Correspondence: chris.brown@otago.ac.nz

¹Department of Biochemistry, Genetics Otago, University of Otago, Dunedin, New Zealand

²Genetics Otago, University of Otago, Dunedin, New Zealand

human 3'UTRs, but can be used for any species sequences, alignments, or any part of the mRNA.

Implementation

The analysis has three phases: 1. accepting user input, 2. analysing the sequence(s), and 3. interactive visualization of the results (Figure 1). The processes to identify and visualise the regulatory elements for any selected gene or given sequence(s) is done in parallel for speed. Input can be the name of a human gene (e.g. TNF) in which case the standard TargetScan/UCSC vertebrate alignment will be used. However, the user can also input any sequence or alignment. The server is a pure LAMP (Linux, Apache, MySQL and Perl) implementation providing speed and stability, using HTML, JavaScript and AJAX to provide seamless user interaction throughout the analysis. SFM has been tested on commonly used web-browsers: Chrome, Firefox, Safari and Explorer 10 or later.

Data analysed

The RNA-Binding Protein DataBase (RBPDB) contains a collection of experimentally verified RNA binding sites, manually curated from literature. It currently contains binding data on 272 RBPs, but only 69 that have motifs in position frequency matrix (PFM) format most useful for SFM analysis. These PFM can be used to distinguish between good and poor matches for short motifs. The other individual binding site sequences from RBPDB could also be user specified (e.g. CAUY). Other user specified sequences, regular expressions, or matrices can also be used in PatSearch format [22].

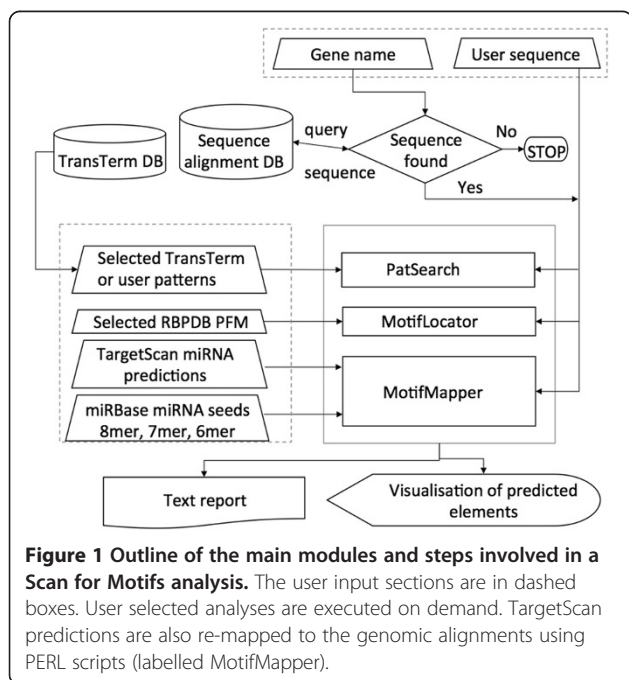


Figure 1 Outline of the main modules and steps involved in a Scan for Motifs analysis. The user input sections are in dashed boxes. User selected analyses are executed on demand. TargetScan predictions are also re-mapped to the genomic alignments using PERL scripts (labelled MotifMapper).

Published miRNA sequences are from miRBase [23]. The mature miRNA sequences were downloaded from miRBase website (file:mature.fa), processed (reverse complemented and 8 leading seed bases extracted) to get a list of 2042 named 8mer seeds and stored in a reference text file. The 6mer seed is the middle 6-bases, and both the two overlapping 7mers are used (7mer-A1, denoted A1 in the output, and 7mer-M8) [8].

The 3'UTR alignments used were obtained from TargetScan (v.6.2) along with the microRNA-binding site related files (miR Family, Predicted Conserved Targets Info, Conserved Family Info) [8]. The 'UTR_Sequences' file holds multiple sequence alignments (MSA) of 23 vertebrate genomes aligned to human, extracted from the USCC human genome (hg18) databases by the TargetScan authors. The human specific sequences were extracted and the positional information for the miR-binding sites provided in "Predicted Conserved Targets Info" file was compared to and updated where needed) against the latest release of hg19 database (from UCSC). A bed format MySQL database table was created to hold the positional information for each of these miR-binding sites.

A custom Perl script was written and used for checking and updating the positional information as above. The program uses sequence similarity between the latest release of hg19 (from UCSC) and the UTR sequences from the TargetScan website. In most of the cases the sequences were 100% identical. For 27 genes the sequences were found to be different in length, the TargetScan prediction data for these were discarded, as they could not be unambiguously assigned to the sequence.

Accepting user input

The user input is of two types, i) query sequence(s) and ii) query element(s). Figure 2 shows the different input options available in SFM web-server.

- i) Query sequence. Option 1 in Figure 2 shows the different types of sequence that is accepted by SFM. It supports input of a standard human gene symbol (i.e. LIN28A) given as source of the query sequence. In such cases relative sequence alignments of 23 vertebrates (including human) will be retrieved from previously processed sequences using the inputted gene symbol and used as query sequence. Alternately, users can input FASTA/multiFASTA/clustalW alignments as well as tabular multiple sequence alignment (MSA) formatted sequences as query sequence. SFM supports assigning reference sequence when the query sequence has more than one sequence. If a human gene symbol was used to get the input sequence, the reference sequence is assigned to be human. In all other cases, the first sequence is considered to be the reference sequence.

Figure 2 The input section of Scan for motifs showing the range of supported regulatory elements and background controls. For a pre-aligned human 3' UTR (e.g. TNF-NM_000594) it defaults to searching for over 60 TransTerm regulatory elements with expectations of E-value ≤ 0.175 by chance in typical human 3' UTR (~1000 nt) (A in Figure) and TargetScan miRNA binding site predictions for ~150 conserved miRNA families (C). In this case the sites for RNA binding proteins with E-values ≤ 1.0 per thousand (B) and miRBase 8mer seeds (D) are also selected.

ii) Query elements. Option 2. A-E in Figure 2 shows the range of query elements expect value controls available in SFM. All the 77 Transterm elements (option 2. A in Figure 2) are associated with an background Expect-value (E-value) frequency of occurrence per thousand bases. These E-values were calculated by first creating a background set by dinucleotide shuffling a non-redundant set of 18,895 human 3'UTR sequences, then searching these with each of the elements. For example an expect value of 0.175 (the default) corresponds to an expectation that each element may appear on average by chance 0.175 times in a typical analysis of one human 3' UTR of 1000 nt. Elements can be automatically selected/deselected by changing the E-value cutoff (shown in the red box in option 2. A in Figure 2.2). Additionally, users can give their own pattern or sequence motif (e.g. AUAGGGU), which will be searched along with the other selected elements against the query sequence(s) using PatSearch.

Similarly, option 2.B-D (Figure 2) shows the elements from RBPDB, TargetScan and miRBase respectively along with the options to limit the hits based on MotifLocator calculated matches using the 69 RBPDB PFM.

The TargetScan elements are available only when a published human gene symbol is used.

Option 2.E (Figure 2) The default behaviour is only to show elements in non-reference sequences if also found in the reference sequence (e.g. human). This can be disabled using this option.

Processing sequences

Upon receiving the input, SFM searches for the query elements using independent parallel processes, where the output from one process is not affected by another process (Figure 1). Irrespective of the input sequence types, all sequences are converted to FASTA format. The patterns from the selected TransTerm elements and user given pattern(s) are used to search the input sequences using PatSearch [22]. The 69 RNA binding protein PFM from RBPDB are used to search the sequences with MotifLocator [24]. The TargetScan miRNA binding sites and their position of occurrences were retrieved from the MySQL database table (see section 2.2.1) by using the input human gene symbol and mapped on the query sequences using PERL scripts labelled MotifMapper in Figure 1. Based on the user given seed length (6, 7 or 8 nucleotides), a list of seed sequences are created from the 2042 seed sequences. As one seed sequence

can be associated with multiple miRNAs in a family, a non-redundant list of seed sequences was made. These sequences were used to search the query sequence(s) using PERL RegEx (regular expressions). Once all the processes are finished, the results from these processes are combined and sent to the visualisation module.

Interactive output

The output is shown on a scrollable alignment with links to further information and the ability to show or hide specific components of the complex results.

Results and discussion

The SFM web-server analyses sequences that may be aligned vertebrate UTRs, or user inputted sequences or alignments (Figure 1). Five types of elements are searched for in these sequences.

- (i) Regulatory elements from the TransTerm database, which includes relevant UTRSite and ARE elements. This provides a curated collection of CREs that function as translational control elements in mRNAs. The computational models (elements) are selected by the user, and/or filtered on empirically determined background frequencies in a shuffled control set. Matches are identified using PatSearch [22].
- (ii) RBP binding sites represented as position frequency matrices (PFM) from the RBPDB [25]. Matches are identified using MotifLocator [24] with a user specified E-value filter.
- (iii) MicroRNA target sites predicted by TargetScan 6.2 [8]. TargetScan was chosen as it is widely used, and predicts sites on vertebrate alignments
- (iv) Human miRNAs 6 to 8 base seed sequences [23] using MotifMapper. This simple prediction is intended to allow visualisation of most of the potential miRNA binding sites, including likely false positives, if the user desires to.
- (v) User defined patterns in PatSearch format [21]. PatSearch allows searches for simple strings, optionally with mismatches insertions and deletions (e.g. GNGNCC), but also more complex elements (e.g. GCG 3...7 GCG, two GCG separated by 3–7 bases) and RNA secondary structures (e.g. p1 = 10...10 4...7 ~ p1, a ten base stem with a loop of 4–7 bases). A full description of the syntax is presented in the help on the SFM server.

On completion of the individual processes, the results are compiled and presented as interactive visualisation (Figure 3). As an example, we use the well-studied tumor necrosis factor alpha (TNF) 3' UTR. TNF is a multifunctional cytokine, it regulates the expression of

other genes in inflammation and other processes and its expression is regulated at main steps [26]. The TNF 3' UTR has been shown to be targeted by both proteins and miRNA [13,27] and is a classic example of an ARE containing mRNA. MicroRNAs that are confirmed to target this UTR in mammals are miR-16 [28], miR-19a [29], miR-125b [30], miR-130 [31], miR-181a [32], miR-301 [31]. Unusually, a miR-369-3p containing RNA-protein complex binds to targets within the ARE and activates or represses translation in the cell cycle [13]. This ARE may also be bound by the RNABP tristetraprolin (TTP) to repress translation [33].

In the SFM analysis using the settings in Figure 2, highlights several types of elements from the TransTerm database (Figure 3, yellow): the AU rich element (ARE) is represented by hits from three overlapping descriptions (Background E-value per thousand bases 0.06, 0.12, 0.12 respectively, Figure 3) [34]; TNF Alpha Stability and Efficiency Element (E-value 0.000008) [35]; and two descriptions of a Polyadenylation Element at the 3' end (E-value 0.03, 0.02). These are all present in a similar position in the alignment across vertebrates, and the 9–12 base core ARE [34] is repeated [34]. The two predicted stability elements in the TNF 3' UTR have been verified experimentally [27,35], and the polyadenylation signal has a clear match to the consensus (AAUAAA). In addition a 15-LOX-DICE element is predicted (E-value 0.01) in the same location in only 5 of 17 species. From the information linked from the small 'i' to the TransTerm entry it can be found that the 15-LOX-DICE is known to have a role in regulating mRNA stability of mRNAs in early erythropoiesis [36]. This may be a false positive, or a novel finding requiring further investigation.

Three predicted overlapping miRNA binding sites are shown (Figure 3, red). Interesting they flank the ARE. Each site links to the family of miRNAs that could bind this seed (e.g. miR181abcd/462) this data is inherited from the TargetScan families and predictions [8]. Included in these predictions are miR-19a, miR-181a, miR-130/miR-301 they have been shown to target these regions in the TNF UTR.

Not predicted with the conservative default SFM parameters are two sites for miR-369-3p within the ARE [13]. These could be shown when 7mer miRBase seeds (miR-369-3p, UAUUAUU) are selected overlapping the ARE. These miR-369-3p sites are also conserved in the alignment. The TargetScan analysis with 153 'broadly conserved' and 'conserved' miRNA families did not predict this site, as miR-369 is poorly conserved [8] so they are not shown in the results from this analysis (Figure 3 red). However, TargetScan does not predict this known site at all (TargetScan webserver) possibly due to the weak AU base pairing within this site.

Such short matches (6mer, 7mer) should be interpreted with caution, as there are over 4000 possible

Authors' contributions

AB designed and developed the software. CMB conceived of the application, supervised it, and tested it. Both authors wrote, read and approved the final manuscript.

Acknowledgements

This work was partially funded by a Human Frontier Science Foundation Research Grant [RGP0031/2009 to Ian Macara, Anne Spang and C.M.B.]; A.B. was a recipient of a University of Otago Postgraduate Scholarship and Publishing Bursary.

Funding

This work was partially funded by a Human Frontier Science Foundation Research Grant [RGP0031 2009 to Ian Macara, Anne Spang and C.M.B.]; A.B. is a recipient of a University of Otago Postgraduate Scholarship.

Received: 10 February 2014 Accepted: 16 May 2014

Published: 8 June 2014

References

- Jacobs GH, Chen A, Stevens SG, Stockwell PA, Black MA, Tate WP, Brown CM: **Transrm: a database to aid the analysis of regulatory sequences in mRNAs.** *Nucleic Acids Res* 2009, **37**(Database issue):D72–D76.
- Chang TH, Huang HY, Hsu JB, Weng SL, Horng JT, Huang HD: **An enhanced computational platform for investigating the roles of regulatory RNA and for identifying functional RNA motifs.** *BMC Bioinforma* 2013, **14**(Suppl 2):S4.
- Szostak E, Gebauer F: **Translational control by 3'-UTR-binding proteins.** *Brief Funct Genomics* 2013, **12**(1):58–65.
- Michalova E, Vojtesek B, Hrstka R: **Impaired pre-mRNA processing and altered architecture of 3' untranslated regions contribute to the development of human disorders.** *Int J Mol Sci* 2013, **14**(8):15681–15694.
- Stevens S, Brown C: **In silico estimation of translation efficiency in human cell lines: potential evidence for widespread translational control.** *PLoS One* 2013, **8**(2):e57625.
- Gruber AR, Fallmann J, Kratochvill F, Kovarik P, Hofacker IL: **AREsite: a database for the comprehensive investigation of AU-rich elements.** *Nucleic Acids Res* 2011, **39**(Database issue):D66–D69.
- Stevens S, Brown C: **Bioinformatic methods to discover cis-regulatory elements in mRNAs.** In *Springer Handbook of Bio-/Neuro-informatics*. Edited by Kasabov N. Heidelberg: Springer; 2014:151–169.
- Garcia DM, Baek D, Shin C, Bell GW, Grimson A, Bartel DP: **Weak seed-pairing stability and high target-site abundance decrease the proficiency of Isy-6 and other microRNAs.** *Nat Struct Mol Biol* 2011, **18**(10):1139–1146.
- Dweep H, Sticht C, Gretz N: **In-silico algorithms for the screening of possible microRNA binding sites and their interactions.** *Curr Genomics* 2013, **14**(2):127–136.
- Naifang S, Minping Q, Minghua D: **Integrative approaches for microRNA target prediction: combining sequence information and the paired mRNA and miRNA expression profiles.** *Curr Bioinform* 2013, **8**(1):37–45.
- Incarnato D, Neri F, Diamanti D, Oliviero S: **MREdictor: a two-step dynamic interaction model that accounts for mRNA accessibility and Pumilio binding accurately predicts microRNA targets.** *Nucleic Acids Res* 2013, **41**(18):8421–8433.
- Ciafre SA, Galardi S: **microRNAs and RNA-binding proteins: a complex network of interactions and reciprocal regulations in cancer.** *RNA Biol* 2013, **10**(6):935–942.
- Vasudevan S, Tong Y, Steitz JA: **Switching from repression to activation: microRNAs can up-regulate translation.** *Science* 2007, **318**(5858):1931–1934.
- Dethoff EA, Chugh J, Mustoe AM, Al-Hashimi HM: **Functional complexity and regulation through RNA dynamics.** *Nature* 2012, **482**(7385):322–330.
- Kedde M, van Kouwenhove M, Zwart W, Oude Vrielink JA, Elkon R, Agami R: **A Pumilio-induced RNA structure switch in p27-3' UTR controls miR-221 and miR-222 accessibility.** *Nat Cell Biol* 2010, **12**(10):1014–1020.
- Wu X, Chesoni S, Rondeau G, Tempesta C, Patel R, Charles S, Dagainawala N, Zucconi BE, Kishor A, Xu G, Shi Y, Li ML, Irizarry-Barreto P, Welsh J, Wilson GM, Brewer G: **Combinatorial mRNA binding by AUF1 and Argonaute 2 controls decay of selected target mRNAs.** *Nucleic Acids Res* 2013, **41**(4):2644–2658.
- Jiang P, Singh M, Collier HA: **Computational assessment of the cooperativity between RNA binding proteins and MicroRNAs in transcript decay.** *PLoS Comput Biol* 2013, **9**(5):e1003075.
- Zhang C, Lee KY, Swanson MS, Darnell RB: **Prediction of clustered RNA-binding protein motif sites in the mammalian genome.** *Nucleic Acids Res* 2013, **41**(14):6793–6807.
- Bryan K, Terrile M, Bray IM, Domingo-Fernandez R, Watters KM, Koster J, Versteeg R, Stallings RL: **Discovery and visualization of miRNA-mRNA functional modules within integrated data using bicluster analysis.** *Nucleic Acids Res* 2014, **42**(3):e17.
- Dassi E, Malossini A, Re A, Mazza T, Tebaldi T, Caputi L, Quattrone A: **AURA: atlas of UTR regulatory activity.** *Bioinformatics* 2012, **28**(1):142–144.
- Grillo G, Turi A, Licciulli F, Mignone F, Liuni S, Banfi S, Gennarino VA, Horner DS, Pavesi G, Picardi E, Pesole G: **UTRdb and UTRsite (RELEASE 2010): a collection of sequences and regulatory motifs of the untranslated regions of eukaryotic mRNAs.** *Nucleic Acids Res* 2010, **38**(Database issue):D75–D80.
- Grillo G, Licciulli F, Liuni S, Sbisà E, Pesole G: **PatSearch: a program for the detection of patterns and structural motifs in nucleotide sequences.** *Nucleic Acids Res* 2003, **31**(13):3608–3612.
- Kozomara A, Griffiths-Jones S: **miRBase: integrating microRNA annotation and deep-sequencing data.** *Nucleic Acids Res* 2011, **39**(Database issue):D152–D157.
- Claeys M, Storms V, Sun H, Michoel T, Marchal K: **MotifSuite: workflow for probabilistic motif detection and assessment.** *Bioinformatics* 2012, **28**(14):1931–1932.
- Cook KB, Kazan H, Zuberi K, Morris Q, Hughes TR: **RBPDB: a database of RNA-binding specificities.** *Nucleic Acids Res* 2011, **39**(Database issue):D301–D308.
- Giambelluca M, Rollet-Labelle E, Bertheau-Mailhot G, Lafamme C: **Post-transcriptional regulation of tumour necrosis factor alpha biosynthesis: Relevance to the pathophysiology of rheumatoid arthritis.** *OA Inflammation* 2013, **1**(1):3.
- Shi JX, Su X, Xu J, Zhang WY, Shi Y: **HuR post-transcriptionally regulates TNF-alpha-induced IL-6 expression in human pulmonary microvascular endothelial cells mainly via tristetraprolin.** *Respir Physiol Neurobiol* 2012, **181**(2):154–161.
- Jing Q, Huang S, Guth S, Zarubin T, Motoyama A, Chen J, Di Padova F, Lin SC, Gram H, Han J: **Involvement of microRNA in AU-rich element-mediated mRNA instability.** *Cell* 2005, **120**(5):623–634.
- Liu M, Wang Z, Yang S, Zhang W, He S, Hu C, Zhu H, Quan L, Bai J, Xu N: **TNF-alpha is a novel target of miR-19a.** *Int J Oncol* 2011, **38**(4):1013–1022.
- Tili E, Michaille JJ, Cimino A, Costinean S, Dumitru CD, Adair B, Fabbri M, Alder H, Liu CG, Calin GA, Croce CM: **Modulation of miR-155 and miR-125b levels following lipopolysaccharide/TNF-alpha stimulation and their possible roles in regulating the response to endotoxin shock.** *J Immunol* 2007, **179**(8):5082–5089.
- Bak RO, Mikkelsen JG: **Regulation of cytokines by small RNAs during skin inflammation.** *J Biomed Sci* 2010, **17**:53.
- Li H, Chen X, Guan L, Qi Q, Shu G, Jiang Q, Yuan L, Xi Q, Zhang Y: **MIRNA-181a regulates adipogenesis by targeting tumor necrosis factor-alpha (TNF-alpha) in the porcine model.** *PLoS One* 2013, **8**(10):e71568.
- Qi MY, Wang ZZ, Zhang Z, Shao Q, Zeng A, Li XQ, Li WQ, Wang C, Tian FJ, Li Q, Zou J, Qin YW, Brewer G, Huang S, Jing Q: **AU-rich-element-dependent translation repression requires the cooperation of tristetraprolin and RCK/P54.** *Mol Cell Biol* 2012, **32**(5):913–928.
- Halees AS, El-Badrawi R, Khabar KS: **ARED Organism: expansion of ARED reveals AU-rich element cluster variations between human and mouse.** *Nucleic Acids Res* 2008, **36**(Database issue):D137–D140.
- Hel Z, Di Marco S, Radzioch D: **Characterization of the RNA binding proteins forming complexes with a novel putative regulatory region in the 3'-UTR of TNF-alpha mRNA.** *Nucleic Acids Res* 1998, **26**(11):2803–2812.
- Thiele BJ, Berger M, Huth A, Reimann I, Schwarz K, Thiele H: **Tissue-specific translational regulation of alternative rabbit 15-lipoxygenase mRNAs differing in their 3'-untranslated regions.** *Nucleic Acids Res* 1999, **27**(8):1828–1836.

doi:10.1186/1471-2105-15-174

Cite this article as: Biswas and Brown: Scan for Motifs: a webserver for the analysis of post-transcriptional regulatory elements in the 3' untranslated regions (3' UTRs) of mRNAs. *BMC Bioinformatics* 2014 **15**:174.