

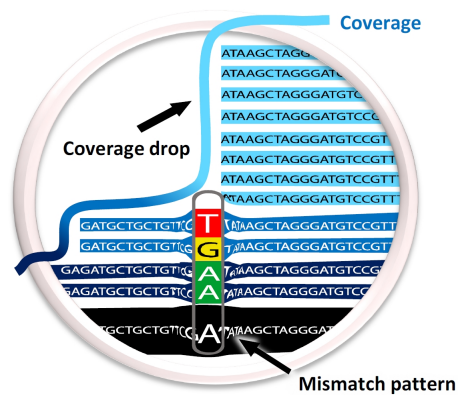
CoverageAnalyzer

Visual Inspection & Automatic Screening of Deep-Seq Data
for Detection of RNA Modifications

v. 1.0

02.2016

Manual



released by

Ralf Hauenschild

from
Obernburg am Main

on 1st February 2016

Prof. Dr. Mark Helm

Institute of Pharmacy and Biochemistry

Johannes Gutenberg-University, Mainz

1 Installation

CoverageAnalyzer is available for Windows, Linux and MacOSX 64 bit systems. The software has several dependencies, which are installed during setup. 2 GB of disk space should be available for installation. Reserve additional space for intended analyses.

1. Download CoverageAnalyzer at <https://sourceforge.net/projects/coverageanalyzer/files/>.
2. Copy to a desired installation folder with full permissions (read/write/execute) and short absolute path (< 30 characters) to avoid installation problems.
3. It is recommended to temporarily disable firewalls/antivirus software or create exceptions.
4. Ensure internet connection and run setup as *administrator*:
 - Windows: Double-click sfx.exe archive to extract and install contents.
 - Linux/MacOSX: Unzip .zip archive and run Setup.sh.
5. Read and accept licence.
6. Wait for download and installation
7. Install Java RE 7+, if prompted.
8. Launch. (Windows: *CoverageAnalyzer.exe*; Linux/MacOSX: *CoverageAnalyzer.sh*)

Troubleshooting:

When setup fails due to accidental launch without administrator rights, kill untermiated processes using Task Manager, such as Miniconda-Setup, vanish install folder and restart setup as administrator. Java RE is a mandatory requirement. When encountering launching problems, make sure Java RE is installed and available via path/environment variable. If Java is not found, try to launch CoverageAnalyzer by `'.../yourjavafolder/java -jar .../yourfolder/CoverageAnalyzer.jar'` or right-click .jar file and specify JRE to open with.

Dependencies:

- Installed by user: Java RE
- Installed automatically: Python (Miniconda), Matplotlib, Numpy, Scipy, Homebrew (MacOSX only), SAMtools

2 Usage

Launch CoverageAnalyzer in any of the following ways:

- Windows
 - Double-click CoverageAnalyzer.exe
 - CMD `'java -jar ...yourpath/CoverageAnalyzer.jar'`
 - Right-click `'...yourpath/CoverageAnalyzer.jar'` and open with JRE
- Linux/MacOSX
 - Terminal: `'sh ...yourpath/CoverageAnalyzer.sh'`
 - Right-click and run/execute `'...yourpath/CoverageAnalyzer.sh'`
 - Terminal: `'java -jar ...yourpath/CoverageAnalyzer.jar'`

2.1 Prerequisites & File Formats

Prior to data inspection, the input files for analysis must be deposited in a folder structure such that a mother folder contains an arbitrary number of sample folders named with integer numbers, each equipped with one *SAM* file named ending with '.sam'. Outside the sample number folders, the original FASTA reference of the *SAM* files must be located named ending with '.fasta'. *BAM* files can be reconverted to *SAM* format using *SAMtools* [3] (see documentation at <http://samtools.sourceforge.net/samtools.shtml>).

2.2 Profiling & Analysis

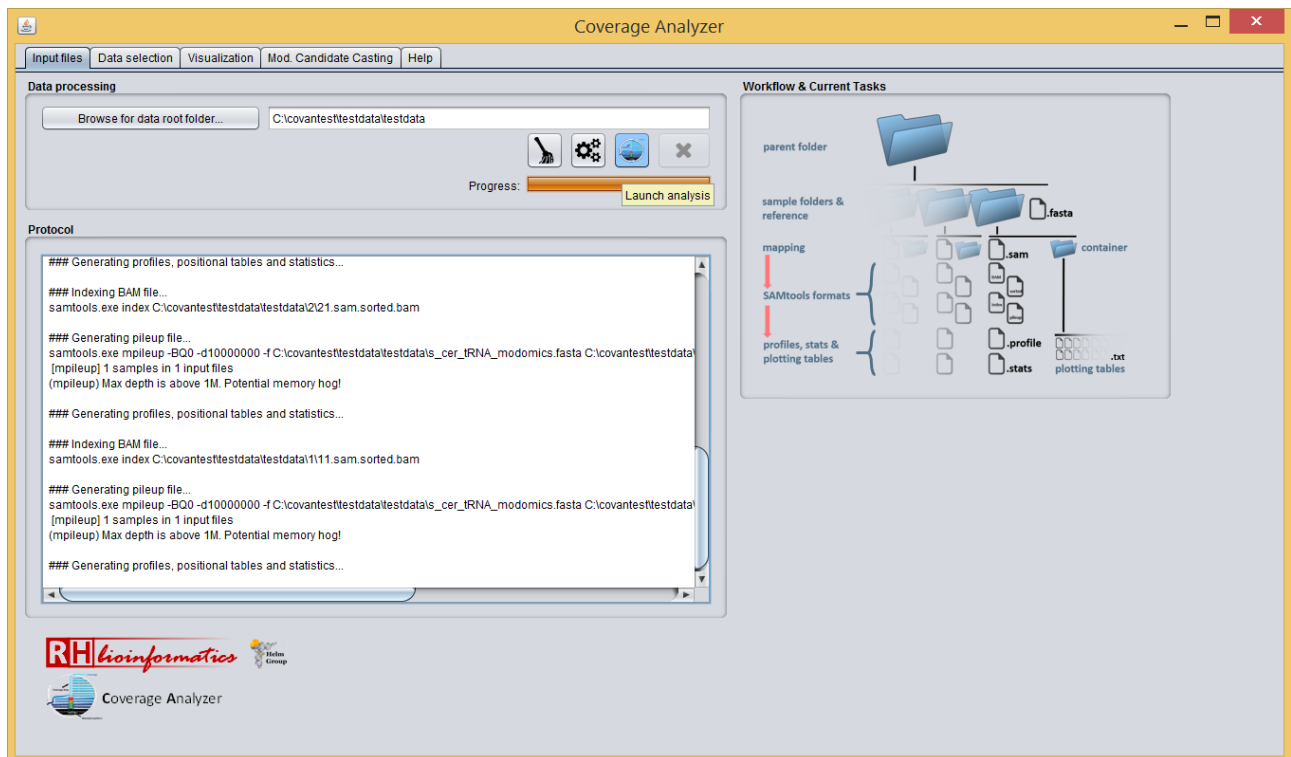
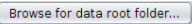






Figure 1: *Input* tab.

CoverageAnalyzer's graphical user interface is fully provided with mouse-hover tool tips. Every analysis starts on the *Input files* tab (Figure 1).

1. Use the browse button  to specify the path of the analysis mother folder containing the sample number folders (choose select/open). Alternatively, a path can be pasted into the adjacent text field. The displayed path should be of the scheme '.../path/motherfolder'. The program will remember the last directory you visited.
2. Optionally, you can vanish the mother folder and subfolders from all data generated previously by CoverageAnalyzer using the  button.
3. Generation of *Profile* format can be performed using the profiling function . Computation time depends on coverage depth in the *SAM*. The current task running is highlighted in the scheme on the right.
4. Pressing of the abortion button  will cancel profiling.
5. When profiling is done, the analysis button  is enabled. Press it and proceed on the *Data selection* tab (Figure 2).

2.3 Table Inspection

The *Data selection* tab allows the user to assess the analyzed mapping profiles with respect to multiple parameters. Next to *sample IDs*, *file paths* and *reference segment* names, also *reference lengths*, *sequences*, *maximum coverage* and the total number of *mapped reads* can be chosen as sorting criteria. A more modification specific evaluation can be done *via* sorting by number of significant arrest sites (coverage ≥ 20 , arrest rate ≥ 0.1), high mismatch rate *m* positions (coverage ≥ 20 , $m \geq 0.1$) and heterogeneous mismatch sites (coverage ≥ 20 , 2nd-most frequent mismatch component ≥ 0.1). Additional requirements, exclusions and cutoffs can be set in the *Filters* section also featuring a restore option for the original table.

Reference s...	Profile path	Reference leng...	Reference seq...	Maximum cove...	# Arrest sites	# Mismatch sites	Hetero-mismat...	# Mapped reads	Sample ID
15IRNAHisjGT...	C:\covantestte...	76	GGCCATCTTA...	9115	24	19	0	9482	2
15IRNAHisjGT...	C:\covantestte...	76	GGCCATCTTA...	9210	16	8	0	8827	1
16IRNAIleIAAT...	C:\covantestte...	77	GGTCTCTTGG...	11207	30	23	1	11572	3
16IRNAIleIAAT...	C:\covantestte...	77	GGTCTCTTGG...	14771	34	23	1	15244	2
16IRNAIleIAAT...	C:\covantestte...	77	GGTCTCTTGG...	14012	20	17	1	14654	3
17IRNAIleIAT...	C:\covantestte...	76	GCTCGGTAG...	2539	10	10	1	2709	3
17IRNAIleIAT...	C:\covantestte...	76	GCTCGGTAG...	4268	10	11	2	4467	2
17IRNAIleIAT...	C:\covantestte...	76	GCTCGGTAG...	4429	8	7	2	4355	1
18IRNAIleuTA...	C:\covantestte...	85	GGGAGITGG...	1520	13	8	0	1785	3
18IRNAIleuTA...	C:\covantestte...	85	GGGAGITGG...	3265	14	8	0	3587	2
18IRNAIleuTA...	C:\covantestte...	85	GGGAGITGG...	5371	16	4	1	5679	1
19IRNAIleuC...	C:\covantestte...	85	GGTGTITGG...	12958	34	28	1	13506	3
19IRNAIleuC...	C:\covantestte...	85	GGTGTITGG...	20063	42	33	0	21207	2
19IRNAIleuC...	C:\covantestte...	85	GGTGTITGG...	23248	18	11	1	25642	1
19IRNAIleuAGC...	C:\covantestte...	76	GGGCGGTG...	10392	28	24	1	10763	3
19IRNAIleuAGC...	C:\covantestte...	76	GGGCGGTG...	13358	30	22	1	13831	2
19IRNAIleuAGC...	C:\covantestte...	76	GGGCGGTG...	10057	17	11	1	10464	1
20IRNAIleuTA...	C:\covantestte...	87	GGAGGTTGG...	2496	21	16	0	3285	3
20IRNAIleuTA...	C:\covantestte...	87	GGAGGTTGG...	5531	19	22	0	6420	2
20IRNAIleuTA...	C:\covantestte...	87	GGAGGTTGG...	10619	21	9	0	11604	1
21IRNAIlysCT...	C:\covantestte...	76	GCCTGTGG...	4620	27	12	1	5129	3
21IRNAIlysCT...	C:\covantestte...	76	GCCTGTGG...	8727	31	16	1	9301	2
21IRNAIlysCT...	C:\covantestte...	76	GCCTGTGG...	21153	15	5	1	21782	1
22IRNAIlysIT...	C:\covantestte...	76	TCCTGTAG...	4043	20	14	1	4307	3
22IRNAIlysIT...	C:\covantestte...	76	XCCTGTAG...	6059	25	16	1	6373	2
22IRNAIlysIT...	C:\covantestte...	76	TCCTGTAG...	7897	15	7	1	8018	1
23IRNAIleuCA...	C:\covantestte...	76	GCTTCAGTAG...	3468	15	11	0	3377	3
23IRNAIleuCA...	C:\covantestte...	76	GCTTCAGTAG...	5324	12	8	0	5562	2
23IRNAIleuCA...	C:\covantestte...	76	GCTTCAGTAG...	6227	10	7	0	6485	1
24IRNAIleuG...	C:\covantestte...	76	GCGGATTTAG...	15587	19	9	1	15932	3
24IRNAIleuG...	C:\covantestte...	76	GCGGATTTAG...	27100	18	9	1	27532	2
24IRNAIleuG...	C:\covantestte...	76	GCGGATTTAG...	28446	14	8	1	28786	1
25IRNAIleuG...	C:\covantestte...	76	GCGGACTTAG...	4455	13	9	1	4703	3
25IRNAIleuG...	C:\covantestte...	76	GCGGACTTAG...	6863	15	9	1	7150	2
25IRNAIleuG...	C:\covantestte...	76	GCGGACTTAG...	7036	10	6	1	7315	1
26IRNAIProITG...	C:\covantestte...	75	GGGCGGTG...	10468	17	27	1	11034	3
26IRNAIProITG...	C:\covantestte...	75	GGGCGGTG...	18506	17	31	1	19130	2
26IRNAIProITG...	C:\covantestte...	75	GGGCGGTG...	44691	13	22	1	42367	1
27IRNAISerjC...	C:\covantestte...	85	GGCAGTATGG...	1038	10	4	0	1247	3
27IRNAISerjC...	C:\covantestte...	85	GGCAGTATGG...	4120	11	5	0	4520	2



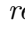

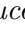
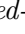
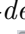
Figure 2: *Data selection* tab.

2.4 Visual Inspection

Click the button in order to launch visualization of a selected alignment. A first plot with standard settings will show up in the *Visualization* tab (Figure 3).

2.4.1 Preferences & Plot Control

The number of plots displayed on one page is changeable to 1 or 2 using the spinner. All specifications (parameters, sequence range, plot details) apply to the plot currently selected in plot selection drop-down menu . The delete button will remove the currently selected plot and resize the remaining ones according to . The button deletes all plots at once. This option is also useful, if erroneous program behavior is observed. Plots can be resized by dragging the main frame of the program window or the right-side divider. Pushing the button or reselecting the plot parameter in the drop-down menu will redraw the selected graph with all changes in settings are realized. If plots appear in wrong aspect ratio after repeated adding or removal, press in order to refresh all plots (this can take a few seconds depending on the number plots and their contents). Preferences can be specified for reportable bases A G T C, as well as for the visual range *r* of the environment adaptive parameter Context

Sensitive Arrest rate, CSA, the fold change of a position i 's arrest rate A w.r.t. median arrest rate within visual range: $CSA^r(i) = \frac{A_i}{\text{median}(A_{i-r}, \dots, A_{i-1}, A_{i+1}, \dots, A_{i+r})}$. In order to display mismatch patterns only at one single sequence position of interest, enable the *center only* checkbox. If needed, the displayed sequence positions can be shifted back or forth using the *pos. shift* spinner. A toggle button optionally shows and hides the arrest rate curve. The sequence interval of interest can be specified in the start and end spinners or by dragging or shifting the range slider. Note that sequence information is plotted up to 100 bases length only, also depending on image width. Zoom in, using the zoom functionality   or the start and end spinners. If longer RNA references are plotted, the software switches to a *reduced-detail mode*, showing coverage only. Display of plot legend , arrest rate curve , reference sequence  and grid lines  as well as additional PDF image generation on the hard drive  can be toggled optionally.

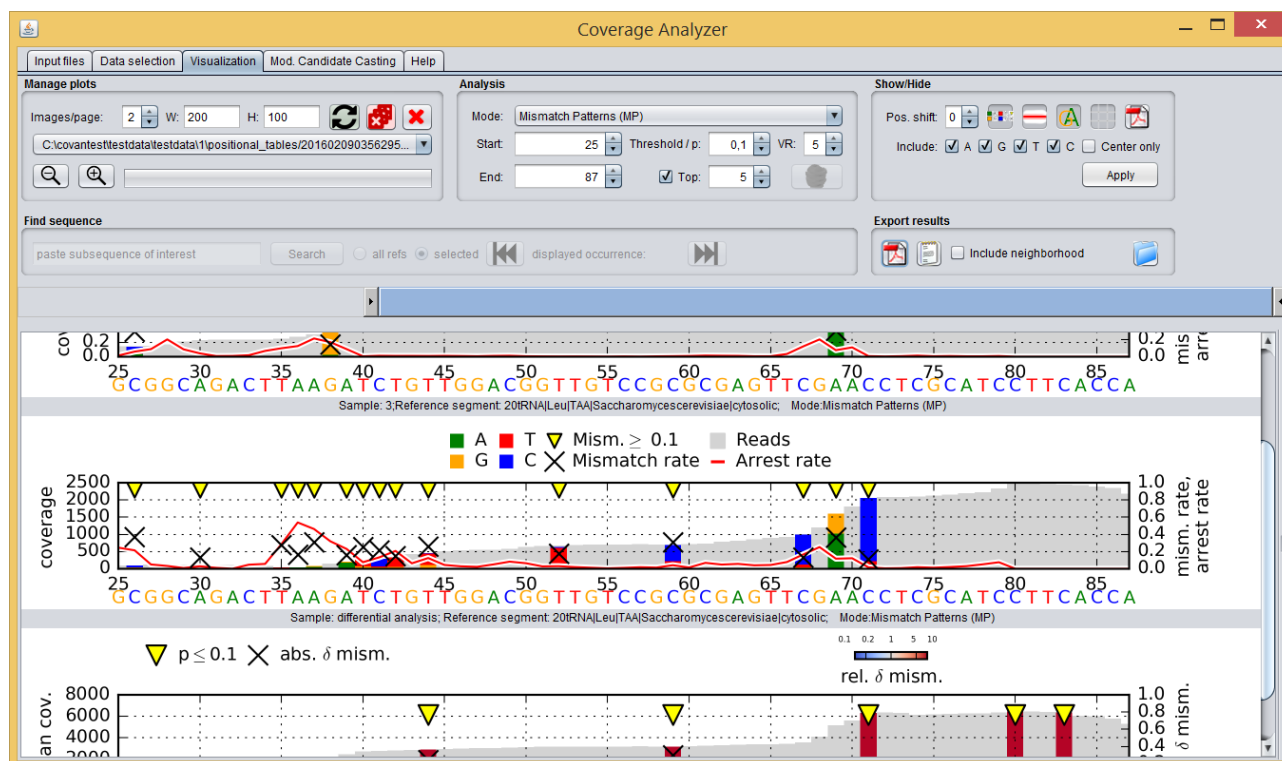



Figure 3: *Visualization* tab. Profiles of two different experimental conditions for the same RNA species were plotted. In addition, a differential plot was generated.




2.4.2 Differential Analysis

Whether for assessment of chemical treatment effects on RT signatures or for comparison of native samples of distinct origins, the differential analysis function is a powerful tool highlighting qualitatively diverging sequence positions. Press button  in order to compare the currently displayed references (e.g. Figure 3). The function is active for equal references only. The readouts are absolute differences indicated by the black markers and relative differences indicated by the colors. Red indicates a stronger signal in the lower plot compared to the upper plot and blue a weaker signal in the lower plot, respectively. Fold changes outside of the .1 to 10 are colored the same as the margin values. If the user wants less events highlighted, but is unsure which pvalue to choose, the usage of for instance a top 5 limit (spinner at upper center) is recommended.

2.4.3 Sequence Search



This feature will be activated in future releases.

2.5 Candidate Screening




The *Mod. Candidate Casting* tab (Figure 4) provides a formula editor for generation of screening filters of arbitrary complexity. Samples can be screened in independent or differential mode, specifically or in batch mode (all). Filter modes control whether screening conditions are applied in require-all mode (and-connected conditions), alternative mode (fulfilling one of the conditions is sufficient for the candidate) or mutual exclusive mode (xor-connected conditions). Alternatively, the user can define a custom formula using parentheses and arbitrarily chosen conditions combined by the boolean buttons   . The intelligent recognition mechanism helps you to set valid brackets (bracket button) around selected pairs of boolean statements by automatic selection snapping. Formulas can be saved as presets and loaded in later sessions. Filter parameters for a candidate position i include coverage c , 3' coverage, arrest rate A , mismatch rate M , M/A ratio, Context Sensitive Arrest rate CSA for specific visual range r and an auxiliary feature termed diversity. *Diversity* d of a single profile position is parameter for experienced users only. Though this feature can be used for other modifications too, its usage is recommended for m¹A screening only. It is represented as a 5-bit binary code, e.g. 10101 with the i^{th} bit set to 1, if the i^{th} of properties

$\left\{ \left[(M \geq 0.1 \wedge M - \max(M^G, M^T, M^C) \geq 0.1) ; [M \geq 0.2] ; [A \geq 0.2] ; [CSA \geq 2] ; [10 \geq M/A \geq 0.1] \right] \right\}$ is fulfilled by the instance, where M is the mismatch rate as sum of mismatch components M^G, M^T, M^C .




Correspondingly, $d := \sum_{i=0}^4 f \cdot 2^i$, where $f = 1$, if condition $5 - i$ is fulfilled and $d \in [0, 31]$ is the decimal representation. Thus, a higher non-m¹A diversity score mimics higher similarity to m¹A. As a quick-setter of this arbitrary set of conditions, the bit code serves two purposes, providing a rough impression of pronounced signature parameters at first glance, while uprating particularly such m¹A candidates meeting the most characteristic features of m¹A, *i.e.* heterogeneous mismatch composition as described in [2], by the implied exponential weighting.

Specific base types can included or discarded from screening report. The *ignore* function allows hiding of certain positions specified by comma-separated integers, e.g. known modification sites out of current interest that would increase the number of reported candidates. Casting is started by the search button . Screening results in *Candidate* format can be accessed *via* the folder button . They are located in the sample number folders chosen for screening. *Candidate* files from differential screening are named according to the compared samples.

2.6 Serial Plotting

Positions of interest for a serial plotting can be provided in *Candidates* format. *Zoom* specifies the number of bases left and right of a candidate to be included in the plot. Having selected the *Candidates* file through the *Batch* button . The plotting series is launched *via* clicking the *Serial plot* button  and the timecode-labeled plots can be inspected in the result folder  (subfolder 'positional_tables' in sample number folder). For every timecode, the reference sequence ID can be obtained from the corresponding '_name' file.

2.7 Export

Positional information on modification candidates and other reported events in the selected plot can be exported *via* the export button . The checkbox *include adjacent positions* controls supplementation of reported positions with the adjacent ones. Button  allows us to jump to the folder containing the saved tab-separated result text file. Button  can be clicked in order to create a PDF file containing all currently displayed plots. For high resolution vector graphics consider the auto-generated PDFs in the plot folder. From these, either JPG snapshots can be generated. Alternatively, vector graphics (EPS recommended) *via* conversion by *e.g.* Inkscape [1] freeware for a flexible graphical post-editing are recommended.

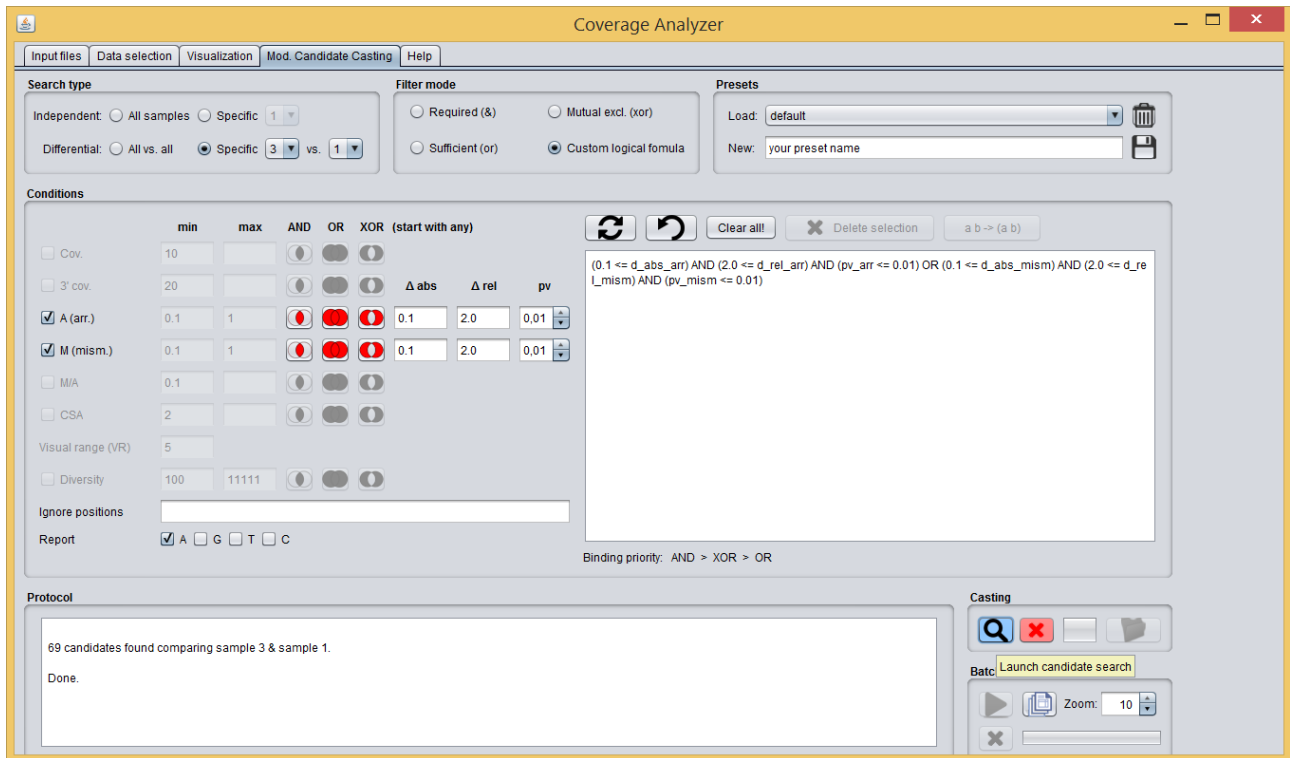



Figure 4: *Mod. candidate casting* tab. A formula for differential analysis of sample 3 vs. sample 1 was generated. By pressing the search button, 69 positions were written to a *Candidates* file on the hard drive.

3 Troubleshooting & Support

Users are encouraged to report any unexpected software behavior to ralf.hauenschild@uni-mainz.de. A common cause of error is invalid formatting of the *Profile* input file. In such case, verification of correct specification and completeness of the tab-separated data should help. In one case, a block of invalid characters turned up at an arbitrary positional line of a reference segment in *Profile*, leading to abortion of the loading process in the *Analysis* step. Reproducibility of this error is not tested yet, but manual correction cured the problem and allowed further study of the data. Other bugs may be resolved by vanishing , followed by restart of the *Analysis* routine (2.2).

4 References

- [1] Inkscape Community. Inkscape.
- [2] Ralf Hauenschild, Lyudmil Tserovski, Katharina Schmid, Kathrin Thüring, Marie-Luise Winz, Sunny Sharma, Karl-Dieter Entian, Ludivine Wacheul, Denis L. J. Lafontaine, James Anderson, Juan Alfonzo, Andreas Hildebrandt, Andres Jäschke, Yuri Motorin, and Mark Helm. The reverse transcription signature of n-1-methyladenosine in rna-seq is sequence dependent. *Nucleic acids research*, 43(20):9950–9964, 2015.
- [3] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, Richard Durbin, et al. The sequence alignment/map format and samtools. *Bioinformatics*, 25(16):2078–2079, 2009.