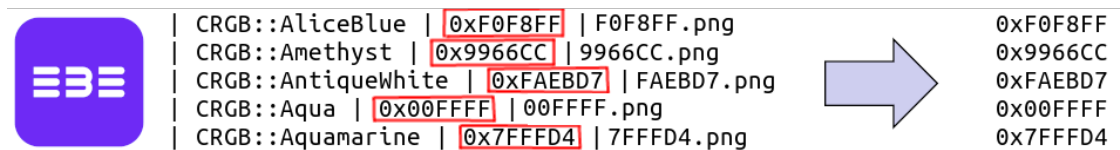


Ebe – A tool for automated file editing using genetic programming

Marek Sedláček*



Abstract

File editing is a big part of today’s work for many people, but not everyone has programming skills or deep knowledge of editing tools to make their editing efficient and quick. This is exactly what Ebe is trying to solve. Ebe takes snippets of file edits done by the user and using genetic programming it finds the correct algorithm to transform the whole file or even multiple files into desired output. Ebe is currently in early version 0.3, but despite that Ebe already achieves some notable results and already contains some additional features for more skilled users to get the most out of it. Ebe is not only for non-programmers, since it can find some edits on its own within seconds, it is a great alternative to handwriting a script for such edits. Even though machine learning is the current hot topic, Ebe uses the approach of evolution – genetic programming – to find the solution, which makes Ebe quite a unique tool and this approach brings in some advantages such as low computational requirements and no need for internet communication with a cloud.

Keywords: Compiler — Interpreter — Genetic programming

Supplementary Material: [Downloadable Code](#)

*xsedla1b@vutbr.cz, Faculty of Information Technology, Brno University of Technology

1. Introduction

Almost everyone who does their work on computer has to work with files and edit them in some ways. It might be as simple as deleting a few lines or as complicated as deleting specific values and moving data or even whole columns around. For programmers, when it comes to those more difficult edits, this might usually mean writing a script to do these edits, but what about the others? Physicist that wants to transform a data set, linguist that want to remove unwanted entries and just overall anyone. Is editing by hand really the best solution? Or forcing these people to learn to program just to do a few edits like this?

Doing small edits is a perfect job for a file editor and hands-on approach, but when editing a file with

thousands of lines, this is no longer a viable option and some other tool needs to be employed. In many cases such “tool” is a program in some programming language, but writing this code requires programming skills and, depending on the skill, notable time to write it and tests it. So a tool that claims to do this, should not only “get the job done”, but also should not take a long time to do so and require deep knowledge from the user, otherwise working with the tool will take longer than the actual work for which the file is edited.

Nowadays there are many methods and tools for automated specialized file editing and file transformation. File types and their formats vary a lot and thus only a one tool cannot handle all these types, but rather specializes on a certain file type and format. But there

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31 still exist tools, which can handle almost all possible
 32 formats, but those tools then require the user to put in
 33 a lot of effort to make them correctly work for all the
 34 possible inputs and cases. Such example can be the
 35 AWK programming language, which was designed in
 36 the late 1980s by Alfred Aho, Peter Weinberger and
 37 Brian Kernighan for the sole purpose of file editing [1].
 38 Since AWK is used till this day and is a built-in tool
 39 for many operating systems [2], this only suggests that
 40 the need has not gone away, but in fact may have even
 41 gone up, with the rise of high computing and large
 42 storage capabilities for computers.

43 Another popular editing approach is the use of gen-
 44 eral purpose programming language, such as Python
 45 3, Perl or even Bash. This approach allows to work
 46 with almost any file format, but requires additional
 47 work to adjust to it and deep enough knowledge of the
 48 language and the file that is being edited.

49 Ebe is trying to tackle all of these problems, it
 50 requires the user to only know how to run it and offers
 51 a quite fast editing (compilation and interpretation)
 52 times compare to writing and running scripts by hand
 53 or even doing all the editing by hand.

54 Currently Ebe is not a perfect do-it-all tool, but
 55 already offers a quite powerful ability to do many edits
 56 in a short time with the option that it will find these
 57 edits for the user. And unlike many other tools, offers
 58 the option to edits multiple similar files at once. It
 59 also uses data types and thus in combination with user
 60 defined expressions allows for changing values in the
 61 whole file based on an expression. And if someone
 62 does not trust the evolution, they can always write the
 63 editing algorithm themselves and then interpret it using
 64 Ebe.

65 2. Genetic programming for code gener- 66 ation

66 Genetic programming (GP) is technique popularized
 67 by John Koza in the 1990s and it is a process of opti-
 68 mization [3]. When speaking about code generation
 69 with genetic programming, the whole process some-
 70 what imitates evolution in the nature. It starts with
 71 randomly initialized populations of candidate pheno-
 72 types, which in this case would be simple programs.
 73 Just as in the nature, these candidates are crossed with
 74 each other, mutated, scored and then the best ones
 75 are selected into future populations based on the scor-
 76 ing [4]. Scoring of the program and the strength of GP
 77 compare to pure randomness can be seen in figure 1.

78 Scoring is a big part of genetic programming and
 79 there always needs to be a scoring function, called the
 80 “fitness function”. This function guides the evolution,

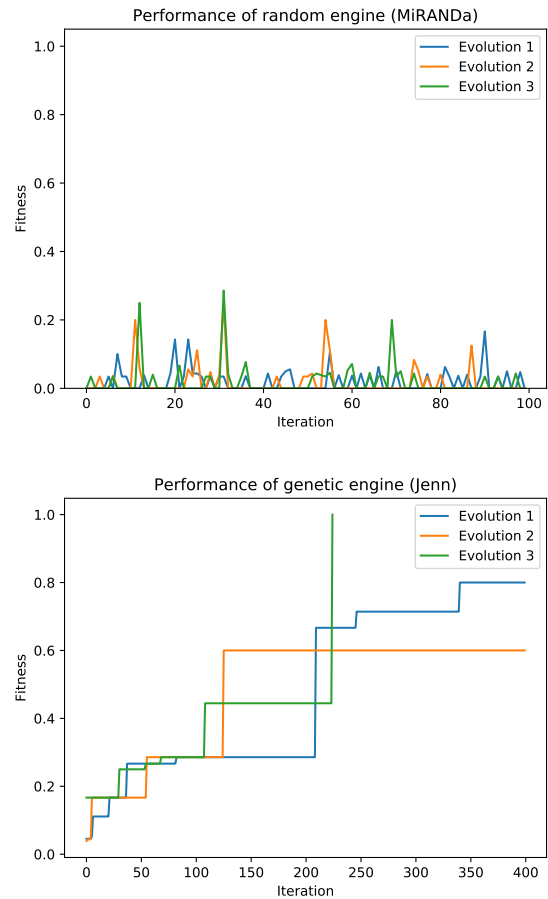


Figure 1. Comparison of 3 iterations of program evaluation using a pure random – MiRANDa engine – program generation (*top graph*) and genetic programming – Jenn engine – approach (*bottom graph*). The GP approach steadily improves the quality of the program, whereas random approach relies only on brute force approach.

by favoring phenotypes with better fitness score [3]. 81
 Without a fitness function, there is no way to guide the 82
 evolution and thus this limits genetic programming to 83
 a limited set of problems, where it can be employed. 84

85 The need for fitness function, somewhat limits the
 86 use for code generation, since often the implemen-
 87 tation of fitness function would result in no longer
 88 needing the GP, since the fitness function would be
 89 the solution itself. An example of this would be when
 90 one was to evolve a sorting algorithm using genetic
 91 programming. To score a solution, the output of the
 92 generated code needs to be compared to some ground
 93 truth, which in this case would be already sorted input
 94 and thus the fitness function could be used for this job
 95 instead of the generated program.

96 But this limitation does not mean that genetic pro-
 97 gramming has no use in this field. On top of generating
 98 code, where fitness function can be provided (such as

99 algorithms for file editing), it can help optimize func-
100 tions (its time or even space complexity), by evolu-
101 ing parts of the code into ways a programmer would
102 not think of [5] or it can also help automatically fix
103 bugs [6]. In other similar fields GP can be used for
104 tasks such as generating computer art [7] ¹, compose
105 music [8] and many other uses.

106 3. Ebe – Edit by Example

107 Ebe is a program (compiler and interpreter) for editing
108 files just from given examples.

109 All the user has to do is create a snippet of file
110 before and after the desired edits (input and output
111 example), give this to Ebe and it will try to find a fitting
112 algorithm (program in *Ebel* language), which would
113 do the desired transformations from input example to
114 the output example (example of Ebe’s usage can be
115 seen in figure 2).

116 Since Ebe is aimed also at people with little to no
117 programming knowledge, Ebe’s philosophy is to not
118 cause exceptions and errors as long as it is not neces-
119 sary. Meaning that, when an incorrect instruction or
120 input is encountered, rather than exiting the execution
121 with error, only a warning is printed and the instruction
122 or input is ignored. Ebe is very verbal about this and
123 will notify about any problems, but this philosophy
124 is quite handy when interpreting multiple files, which
125 might slightly differ (a division by 0 might be encoun-
126 tered or different data type at some position), but one
127 Ebel will work on all of them.

128 Ebe consists of multiple independent modules (see
129 figure 3), where some of them can even be used in-
130 dependently. One of the main modules is *Ebec* (Ebe
131 compiler), which does all the necessary file parsing
132 and guides the evolution, which happens in the *Eben*
133 (Ebe engine). Once the compiler and engine evolve
134 sufficiently fitting program, then this can be outputted
135 or loaded into the *Ebei* (Ebe interpreter), which can
136 then edit any number of files using this algorithm.

137 The time for Ebe to generate a suitable program
138 is non-deterministic and depends on lots of external
139 factors. Some very simple programs can be evolved
140 within tenth of a second and some difficult ones might
141 take a few minutes to be evolved. This might be a
142 significant time to some people, but compare to al-
143 ternative approach with editing by hand or writing
144 an editing script, this might be faster and simpler ap-
145 proach, since it can run in the background, while other
146 work is done.

¹Abstract evolution – a program for generating abstract com-
puter art using genetic programming – <https://github.com/mark-sed/abstract-evolution>.

```
[mark-sed@ebe]$ ls
example.in  example.out  temps.data
[mark-sed@ebe]$ cat temps.data
Prague: 42
Paris: 50
Venice: 60
Munich: 38
[mark-sed@ebe]$ cat example.in
Prague: 42
[mark-sed@ebe]$ cat example.out
Prague: {!( ($ - 32) * 5) / 9 !}
[mark-sed@ebe]$ ebe -in example.in -out ex
Perfectly fitting program found.

Best compiled program has 100% precision
[mark-sed@ebe]$ ebe -i f2c.ebel temps.data
Prague: 5
Paris: 10
Venice: 15
Munich: 3
```

Figure 2. Example of Ebe’s user defined expressions to change values in Fahrenheit to Celsius. This example uses split compilation and interpretation, but it is possible to do both in one Ebe invocation.

4. Ebe’s implementation

147 One of Ebe’s strengths is that it generates Ebel in-
148 structions, which in most cases care only about the
149 word’s position rather than its type or even value. This
150 means that the example provided can only have the
151 same structure as the actual file that needs to be edited.
152 Meaning that one can generate a general editing script
153 by only knowing the structure of the actual input file
154 (see listing 1 and 2). This gives the option to provide
155 someone an editing script without requesting the actual
156 data that needs to be edited, which could otherwise
157 be a problem for security or other reasons. This ap-
158 proach can also be useful in cases, where the output is
159 ambiguous to the input because of the values. 160

Listing 1. Example input file containing actual values.

```
Brno 42 5.0e-8 + True 161
```

Listing 2. Example input file, which is structurally
and type-wise equivalent to an input file in listing 1
and therefore can be used in place of 1 for Ebe.

```
A 0 0.0 + B 162
```

163 On the other hand Ebel language contains con-
164 structs, which allow to match on string value or word’s
165 type. These more advanced scripts can for example be
166 generated using Bee language and bee-hs compiler².
167 This compiler uses a higher level abstraction language,
168 which it then compiles into Ebel and offers this way
169 an alternative scripting language for file editing for
170 programmers.

²Bee-hs – a compiler for Bee language, which is compiled into
Ebel – <https://github.com/mark-sed/bee-hs>.

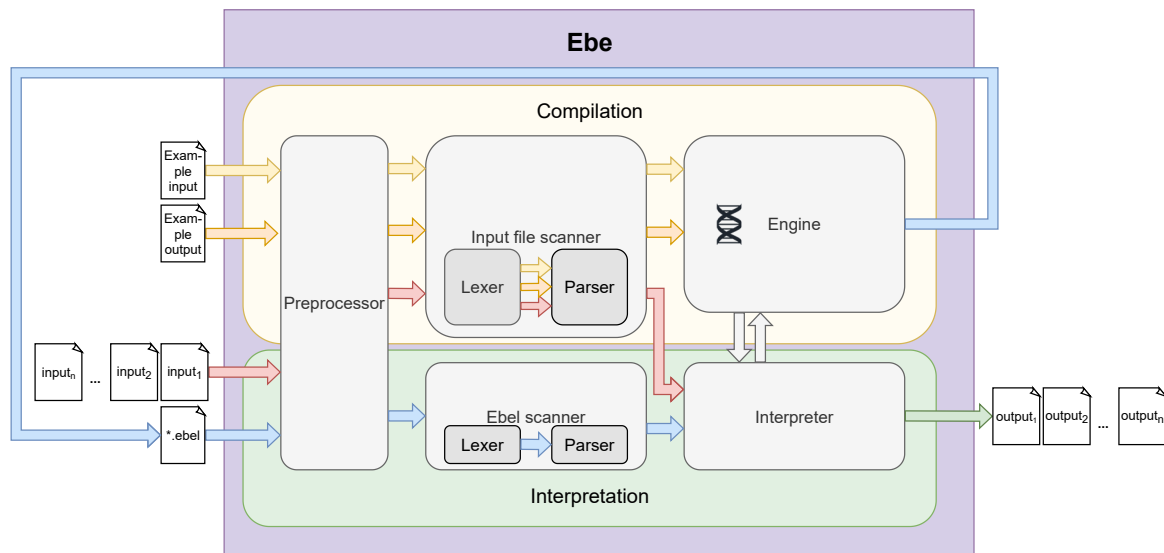


Figure 3. Diagram of Ebe's workflow.

171 The genetic algorithm used in Ebe varies based
 172 on selected engine. Ebe currently offers 3 engines –
 173 the experimental engine *MiRANDa*, which uses only
 174 random walk, engine *Taylor* and engine *Jenn*. Tay-
 175 lor's approach is quite similar to the approach used
 176 in Cartesian genetic programming, where only muta-
 177 tions are used [9] and the starting phenotypes contain
 178 only NOP instructions ("NOP tail"). Jenn is the default
 179 engine and uses genetic programming process quite
 180 similar to the one described by J. Koza [10], where
 181 starting population are random programs of random
 182 size and are during evolution based on set probabilities
 183 crossed over, mutated (where mutation only changes
 184 one instruction for another random one) and then se-
 185 lected into new populations. Fitness is calculated using
 186 one of the string comparison methods (listed below),
 187 which is applied to the example input interpreted using
 188 current phenotype and example output.

189 Ebe is written in C++ and does not use any external
 190 libraries (for general use, otherwise GoogleTest) to al-
 191 low for easier portability and compilation. Ebe is a free
 192 and open-source project and is designed to be modular
 193 and extensible by others. It uses Flex and Bison for
 194 lexer and parser generation and thus allows developers
 195 to easily implement new parser and lexers for new file
 196 formats. The same can be done with engines, where
 197 engines are what powers and does the whole process of
 198 code generation using genetic programming (although
 199 it is possible to use different approach). On top of
 200 this Ebe's evolutionary process can be even controlled
 201 at compile time with multiple command line options
 202 to set attributes such as the fitness function (Leven-
 203 shtein distance, Jaro distance, Jaro-Winkler distance or
 204 "one-to-one" character comparison), population size,
 205 number of generations, number of evolutions or even

a timeout based on compilation time or minimum re- 206
 quired output precision. These options are for more 207
 skilled users and more often meant for experimenting 208
 since Ebe contains heuristics, which decide all these 209
 attributes for the user based on the input. 210

5. Ebel – Ebe language 211

Ebel is an imperative, case insensitive, programming 212
 language designed for file editing and to work well 213
 with genetic programming. 214

Ebel is not really meant to be written, but can be 215
 and contains some syntactic sugar to make writing 216
 and editing it more user-friendly. Ebel resembles a 217
 bytecode and was designed in this way to allow for 218
 quick parsing and execution in the interpreter. It is 219
 interpreted over a file, where the Ebel code can be 220
 thought of as a pipeline of instructions through which 221
 the file objects (lexemes) go and get modified by. 222

Ebel is composed of multiple sections called passes. 223
 A pass defines in which way the input file is read. Each 224
 pass is then composed of instructions which take as an 225
 input objects its parent pass parses (word or line). 226

Listing 3. Example of Ebel code.

```

PASS Words 227
  NOP 228
  DEL 229
  DEL 230
  PASS number Expression 231
    SUB $1, $0, 32 232
    MUL $2, $1, 5 233
    DIV $0, $2, 9 234
    RETURN NOP 235
  PASS derived Expression 236
    RETURN DEL 237
PASS Lines 238
  SWAP 1 239
  LOOP 240
  
```

241 The Ebel code in listing 3 will do the following:

242 1. **PASS Words** - file will be interpreted word by
 243 word and for each line:

244 1.1. **NOP** - 1st object will left as is,
 245 1.2. **DEL** - 2nd object will be deleted,
 246 1.3. **DEL** - 3rd object will be deleted,
 247 1.4. **PASS number Expression** - 4th ob-
 248 ject, if it is a number will be:

249 1.4.1. **SUB** \$1, \$0, 32 - subtract 32 from
 250 its value,
 251 1.4.2. **MUL** \$2, \$1, 5 - multiply the new
 252 result by 5,
 253 1.4.3. **DIV** \$0, \$2, 9 - divide the result by
 254 9 and save it as the new value for the
 255 object,
 256 1.4.4. **RETURN NOP** - end expression and
 257 do not modify the new result.

258 1.5. **PASS derived Expression** - if 4th
 259 object was not a number, then use the fol-
 260 lowing without regarding its type:

261 1.5.1. **RETURN DEL** - delete the object.

262 2. **PASS Lines** - file will be interpreted line by
 263 line and for each line:

264 2.1. **SWAP** 1 - swap current line with the fol-
 265 lowing one,
 266 2.2. **LOOP** - repeat until all lines were processed.

267 As can be seen in listing 3, Ebel contains the means
 268 to carry out computations over a single word (num-
 269 bers, floats and even strings) in the edited file, but
 270 because of the problems with finding correct expres-
 271 sions, where for symbolic regression a large data set
 272 would be needed [11] and even then it could be am-
 273 biguous, this task is left to the user in the form of,
 274 already mentioned, user-defined expressions. If such
 275 expression is defined Ebe treats it as always correct,
 276 but still can evolve the other parts of the file to find a
 277 correct Ebel program.

278 6. Real world Ebe use examples

279 Ebe has not been long enough in a public version, but
 280 there are already some real world cases, where it has
 281 proven to be useful.

282 6.1 Extracting hexadecimal color values from 283 markdown table

284 A documentation for WS2811 LED library³ contained
 285 a table of predefined colors and their hexadecimal val-
 286 ues, which needed to be extracted (see first page teaser

³<https://github.com/FastLED/FastLED/wiki/Pixel-reference>

image for reference). This is a very easy edit for Ebe 287
 and it took only 100 ms to compile the correct Ebel 288
 and interpret it. 289

It required the first line of the table to be put into 290
 the input example file and then edit this by hand in the 291
 output example file, but the overall time is insignificant 292
 to the time it would take to edit this file by hand or 293
 writing a script for this edit. 294

6.2 Adjusting feature indexes in biological 295 data set 296

In this case a desired section of human genome was 297
 extracted from the whole genome in GTF format. But 298
 to display this correctly in a genome browser program, 299
 it needed for the start and end feature index (column 300
 4 and 5) to be moved to a different position. In other 301
 words the value 153350000 was needed to be sub- 302
 tracted from all values in the column 4 and 5. 303

This problem is a perfect case for the use of user 304
 defined expressions in Ebe. Here are the example 305
 input file and example output file (which contains at 306
 the position 4 and 5 the index subtraction): 307

Listing 4. Input example file for ebe (with line breaks 308
 for readability).

```
chr1 hg38_knownGene exon 153357854 308
153357881 0.000000 + . gene_id 309
"ENST00000368738.4"; transcript_id 310
"ENST00000368738.4"; 311
```

Listing 5. Output example file for ebe (with line 312
 breaks for readability).

```
chr1 hg38_knownGene exon {!$-153350000!} 312
{!$-153350000!} 0.000000 + . gene_id 313
"ENST00000368738.4"; transcript_id 314
"ENST00000368738.4"; 315
```

In the output example file (listing 5) the user de- 316
 fined expression are between the {! and !} control 317
 sequences and define the expression to be done over 318
 the value at the expression's position. Since this is 319
 the only edit needed, Ebe compiles and interprets this 320
 within milliseconds. The generated Ebel file can then 321
 be used for any other GTF file (since GTF format is 322
 standardized and does not change structure [12]) and 323
 handles even large 50 MB file (123 705 lines) within 324
 13 seconds time (on Intel® Core™ i5-4690 CPU). 325

7. Conclusions 326

Ebe displays a promise as an editing tool for not only 327
 non-programmers. It requires minimum knowledge 328
 of the tool, finds the algorithm for the user and even 329
 allows to edit multiple similar files at once. 330

331 Although Ebe does not yet contain all the wanted
332 features and optimizations, it still allows for some
333 advanced edits, which it can find in a short time and
334 most importantly it can do large simple edits, which
335 might be lot of times needed more than the complex
336 ones. With the addition of user defined expressions
337 Ebe can also tackle more complex problems.

338 Ebe also greatly showcases the possibilities and
339 power of genetic programming for code generation and
340 as an alternative for deep learning or neural network
341 in some cases.

342 As for the future of Ebe – it is still in development
343 as a free and open source tool for anyone to try it out,
344 with stable releases being published. The goal is to get
345 Ebe to the point, where it can do more complicated
346 edits faster than using other file editing approaches or
347 at least in a similar time frame without any additional
348 help from the user. Graphical interface is also consid-
349 ered for future releases to make Ebe more usable to
350 non-programmers.

351 Since Ebe is open source, anyone is highly encour-
352 aged to try it out, play around with it, change parts of
353 it or even integrate Ebe or parts of it into their own
354 project. It also allows for easy parser extensibility, so
355 new file formats can be easily added into Ebe.

356 Acknowledgements

357 I would like to thank my supervisor prof. Ing. Lukáš
358 Sekanina, Ph.D. for his help and Bc. Klára Ungrová
359 for the design of Ebe's logo.

360 References

- 361 [1] Alfred V. Aho, Brian W. Kernighan, and Peter J.
362 Weinberger. *The AWK Programming Language*.
363 Addison-Wesley Longman Publishing Co., Inc.,
364 USA, 1987.
- 365 [2] Canonical. Ubuntu manpage: Awk - pattern
366 scanning and processing language. <http://manpages.ubuntu.com/manpages/trusty/man1/awk.1posix.html>.
- 369 [3] John R. Koza. *Genetic Programming: On the*
370 *Programming of Computers by Means of Natural*
371 *Selection*. MIT Press, Cambridge, MA, USA,
372 1992.
- 373 [4] Wildor Ferrel and Luis Alfaro. Genetic
374 programming-based code generation for arduino.
375 *International Journal of Advanced Computer Sci-*
376 *ence and Applications*, 11(11), 2020.
- 377 [5] Keith D. Cooper, Philip J. Schielke, and Devika
378 Subramanian. Optimizing for reduced code space

using genetic algorithms. In *Proceedings of the*
ACM SIGPLAN 1999 Workshop on Languages,
Compilers, and Tools for Embedded Systems,
LCES '99, page 1–9, New York, NY, USA,
1999. Association for Computing Machinery.

- [6] Stephanie Forrest, ThanhVu Nguyen, Westley
Weimer, and Claire Le Goues. A genetic pro-
gramming approach to automated software repair.
In *Proceedings of the 11th Annual Conference on*
Genetic and Evolutionary Computation, GECCO
'09, page 947–954, New York, NY, USA, 2009.
Association for Computing Machinery.
- [7] Yuchen Wang and Rong Xie. Pixel-based ap-
proach for generating original and imitating evo-
lutionary art. *Electronics*, 9(8), 2020.
- [8] Dragan Matić. A genetic algorithm for com-
posing music. *Yugoslav Journal of Operations*
Research, 20, 01 2010.
- [9] Julian Miller. *Cartesian Genetic Programming*,
volume 43. 06 2003.
- [10] John R. Koza. Genetic programming - on the
programming of computers by means of natural
selection. In *Complex adaptive systems*, 1993.
- [11] Hitoshi Iba, Ji Feng, and Hossein Izadi Rad. Gp-
svm: Genetic programming-based symbolic regres-
sion using relevance vector machine. In *2018*
IEEE International Conference on Systems, Man,
and Cybernetics (SMC), pages 255–262, 2018.
- [12] GFF/GTF File Format. Definition and supported
options. <http://www.ensembl.org/info/website/upload/gff.html>.