



# Convertor din Markdown in HTML

Git Repository

*Laboratory activity*

Name: Polman Marian, Rus Maria Alina

Group: 30232

Email: marianpolman2000.mp@gmail.com, hexexec@gmail.com

Teaching Assistant: Timotei Molcut  
timoteim98@gmail.com



# Contents

<b>1</b>	<b>Convertor din Markdown in HTML</b>	<b>3</b>
1.1	Introducere . . . . .	3
1.2	Descriere teoretica . . . . .	3
1.2.1	Limbajul Markdown . . . . .	3
1.2.2	Limbajul HTML . . . . .	3
1.2.3	Scopul convertorului . . . . .	4
1.2.4	Functionalitati cheie . . . . .	4
1.3	Implementare . . . . .	4
1.3.1	Metoda de parsare . . . . .	4
1.3.2	Structura de date pentru stocarea conținutului . . . . .	4
1.3.3	Traducerea elementelor Markdown în HTML . . . . .	5
1.3.4	Gestionarea stilizării . . . . .	5
1.3.5	Extensibilitate și configurabilitate . . . . .	5
1.3.6	Arborele de parsare . . . . .	7
1.4	Rezultate . . . . .	9

# Chapter 1

## Convertor din Markdown in HTML

### 1.1 Introducere

Proiectul nostru constă în dezvoltarea unui convertor eficient și flexibil din limbajul Markdown în limbajul HTML. Markdown este un limbaj simplu și ușor de utilizat pentru formatarea textului, în timp ce HTML oferă posibilități mai avansate de structurare și stilizare a conținutului web. Convertorul nostru facilitează procesul de transformare a documentelor scrise în Markdown în pagini HTML bine structurate și vizualmente atractive.

Prin intermediul convertorului nostru, utilizatorii pot crea conținut în Markdown, utilizând sintaxa și convențiile specifice acestui limbaj. Apoi, convertorul preia documentele Markdown și le convertește în cod HTML corespunzător, utilizând etichete și structuri HTML adecvate pentru a reda corect formatarea, stilurile și elementele de structură din documentul original.

Funcționalități cheie ale convertorului nostru includ recunoașterea antetelor de diferite niveluri, inclusiv stilizarea acestora în etichetele corespunzătoare de antet HTML. De asemenea, convertorul detectează textul italic și boldat, îl formatează corespunzător în etichete HTML și menține structura paragrafelor și linii goale pentru o experiență de citire plăcută.

Prin intermediul acestui proiect, ne-am propus să oferim utilizatorilor o modalitate simplă și eficientă de a transforma conținutul scris în Markdown în pagini web HTML, reducând astfel nevoia de a efectua manual conversia și optimizând procesul de publicare a conținutului online.

Sursa limbaj Markdown

### 1.2 Descriere teoretică

#### 1.2.1 Limbajul Markdown

Markdown este un limbaj simplu și intuitiv, folosit pentru a formata textul în mod structurat și ușor de citit. Este bazat pe utilizarea caracterelor speciale și convenții de formatare pentru a indica elemente precum antete, stiluri de text și structura paragrafelor. Markdown permite utilizatorilor să creeze conținut într-un format ușor de scris și de editat.

#### 1.2.2 Limbajul HTML

HTML (HyperText Markup Language) este limbajul standard utilizat pentru crearea și afișarea paginilor web. HTML oferă etichete semantice și structurate care permit definirea ele-

mentelor și stilurilor pentru o prezentare adecvată a conținutului web. Prin utilizarea etichetelor HTML, este posibilă structurarea și stilizarea conținutului în mod flexibil și personalizat.

### 1.2.3 Scopul convertorului

Scopul convertorului nostru este de a facilita procesul de transformare a documentelor scrise în Markdown în pagini HTML bine structurate și vizualmente atractive. Convertorul preia conținutul scris în Markdown și îl traduce în cod HTML corespunzător, astfel încât să fie afișat într-un mod adecvat pe paginile web.

### 1.2.4 Funcționalități cheie

Convertorul nostru recunoaște și convertește următoarele elemente din Markdown în HTML:

Antete: Convertorul detectează antetele de diferite niveluri (de la h1 la h6) și le traduce în etichete HTML corespunzătoare.

Text italic și boldat: Convertorul identifică textul scris între caracterele `"*"` sau `"_"` pentru italic și textul scris între caracterele `"**"` sau `"__"` pentru boldat și îl formatează în etichete HTML corespunzătoare.

Paragrafe: Convertorul menține structura paragrafelor și linii goale pentru o afișare corectă a conținutului.

Alte elemente: Convertorul poate fi extins pentru a recunoaște și converti alte elemente specifice din Markdown, cum ar fi liste sau link-uri.

## 1.3 Implementare

Pentru acest proiect ne-am ajutat de parserul C dat ca exemplu în laborator în special pentru partea de transmitere a textului. Am folosit un union în fisierul `yacc(p.y)` pentru a stoca un element de tip `char*`.

### 1.3.1 Metoda de parsare

Am ales să implementăm un parser Markdown bazat pe expresii regulate pentru a recunoaște și a extrage elementele specifice din conținutul Markdown. Am dezvoltat o serie de reguli de analiză pentru a identifica antetele, stilurile de text și paragrafele în conținutul Markdown și pentru a le transforma în structuri de date interne reprezentând AST-ul.

### 1.3.2 Structura de date pentru stocarea conținutului

Pentru a stoca și a manipula conținutul într-un mod eficient, am implementat o structură de date în formă de arbore pentru reprezentarea AST-ului. Fiecare nod din arbore reprezintă un element specific din Markdown, cum ar fi un antet sau un paragraf, și conține informații relevante, cum ar fi nivelul antetului sau textul paragrafului.

### 1.3.3 Traducerea elementelor Markdown în HTML

Am implementat o funcție de generare a codului HTML, care traversează AST-ul și transformă fiecare nod în cod HTML corespunzător. De exemplu, pentru un antet de nivel 1, generăm eticheta ‘h1’ în HTML, iar pentru stilul de text italic, generăm eticheta ‘em’. Am asigurat că sintaxa și convențiile specifice ale limbajului Markdown sunt respectate în rezultatul HTML generat.

### 1.3.4 Gestionarea stilizării

Pentru gestionarea stilizării textului în HTML, am implementat o funcție care recunoaște caracteristicile specifice ale Markdown pentru italic și boldat și le convertește în etichete HTML adecvate. Astfel, textul scris între asteriscuri sau sublinieri simple (‘\*text\*’ sau ‘\_text\_’) este încapsulat în eticheta ‘em’, iar textul scris între dublu asteriscuri sau duble sublinieri (‘\*\*text\*\*’ sau ‘\_\_text\_\_’) este încapsulat în eticheta ‘strong’.

### 1.3.5 Extensibilitate și configurabilitate

Pentru a permite utilizatorilor să personalizeze convertorul, am adăugat opțiuni de configurare suplimentare. Utilizatorii pot specifica adăugarea de stiluri CSS personalizate sau pot activa/deactiva anumite opțiuni de conversie prin intermediul unui fișier de configurare. Am implementat un sistem flexibil de gestionare a configurării, care permite adaptabilitate și extensibilitate.

```
1 %{
2     #include <stdio.h>
3     #include <stdlib.h>
4     #include <string.h>
5     #include "y.tab.h"
6 }
7
8 %option noyywrap
9 string ([a-zA-Z1-9.?! ])*
10 heading1 ("#")
11 heading2 ("##")
12 heading3 ("###")
13 heading4 ("####")
14 heading5 ("#####")
15 heading6 ("#####")
16 bold ("_")
17 bold2 ("**")
18 italic ("")
19 italic2 ("*")
20 br (" ")
21 NL \n
22 %%
23 {italic}      return ITALIC;
24 {italic2}     return ITALIC;
25 {NL} return NL;
26 {br}+ return BR;
27 {bold}|{bold2} return BOLD;
28 {string} {yyval.string = yytext; return STRING;};
29 {heading1} return HEADING1;
30 {heading2} return HEADING2;
31 {heading3} return HEADING3;
32 {heading4} return HEADING4;
33 {heading5} return HEADING5;
34 {heading6} return HEADING6;
```

## Listing 1.1: p.l

```

35 %{
36 #include <stdio.h>
37 #include <stdlib.h>
38 #include <stdarg.h>
39 #include <string.h>
40 char* removeChar(char* s, char c);
41 %}
42 %union{
43     char* string;
44 }
45 %token <string> STRING
46 %token HEADING1 HEADING2 HEADING3 HEADING4 HEADING5 HEADING6 BR NL BOLD
    ITALIC
47 %%
48 FISIER: INTRO BODYPART OUTRO;
49 BODYPART: BREAK | HEADING | P | BODYPART BREAK | BODYPART HEADING | BODYPART
    P | ITA | BODYPART ITA | BLD | BODYPART BLD;
50 BREAK: BR NL {printf("    <br>\n");};
51 P: STRING NL NL {printf("    <p>%s</p> \n", removeChar($1, '\n'));};
52 ITA: ITALIC STRING ITALIC NL {char* s=removeChar($2, '\n');
53                               printf("    <em>%s</em>\n", removeChar(s, s[
    strlen(s)-1]))};
54 BLD: BOLD STRING BOLD NL {char* s=removeChar($2, '\n');
55                             printf("    <strong>%s</strong>\n", removeChar(s, s[strlen(s)
    -1]))};
56 HEADING: HEADING1 STRING NL {printf("    <h1>%s</h1> \n", removeChar($2, '\n'))
    ;}
57 | HEADING2 STRING NL {printf("    <h2>%s</h2> \n", removeChar($2, '\n'))};
58 | HEADING3 STRING NL {printf("    <h3>%s</h3> \n", removeChar($2, '\n'))};
59 | HEADING4 STRING NL {printf("    <h4>%s</h4> \n", removeChar($2, '\n'))};
60 | HEADING5 STRING NL {printf("    <h5>%s</h5> \n", removeChar($2, '\n'))};
61 | HEADING6 STRING NL {printf("    <h6>%s</h6> \n", removeChar($2, '\n'))};
62 INTRO: {printf("<!DOCTYPE html>\n");
63           printf("<html> \n");
64           printf("<body> \n");
65           };
66 OUTRO: {
67           printf("</body> \n");
68           printf("</html> \n");
69 };
70 %%
71
72 char* removeChar(char* s, char c)
73 {
74
75     int j, n = strlen(s);
76     for (int i = j = 0; i < n; i++)
77         if (s[i] != c)
78             s[j++] = s[i];
79
80     s[j] = '\0';
81     return s;
82 }
83
84
85 #include "lex.yy.c"
86 #include <ctype.h>

```

```

87 int main()
88 {
89     yyparse();
90 }
91 int yyerror() {}

```

Listing 1.2: p.y

### 1.3.6 Arborele de parsare

Utilizarea site-ului ANTLR pentru afișarea arborelui de parsare este o modalitate excelentă de a vizualiza structura internă a conținutului Markdown și modul în care este interpretat de parser. Arborele de parsare oferă o reprezentare grafică a ierarhiei elementelor și relațiilor dintre acestea.

Site-ul ANTLR oferă o interfață ușor de utilizat pentru a încărca fișierul de gramatică și a genera arborele de parsare. După ce arborele este generat, acesta poate fi explorat, extins și restrâns pentru a vizualiza diferitele niveluri ale ierarhiei. Aceasta permite o înțelegere mai clară a modului în care regulile de parsare din gramatică sunt aplicate și cum sunt identificate elementele specifice din conținutul Markdown.

```

92 lexer grammar ExprLexer;
93
94 STRING: [a-zA-Z1-9.?! ]+;
95 HEADING1: '#';
96 HEADING2: '##';
97 HEADING3: '###';
98 HEADING4: '####';
99 HEADING5: '#####';
100 HEADING6: '#####';
101
102 BOLD: '_';
103 BOLD2: '**';
104 ITALIC: '*';
105 BR: ' ';
106 NL: '\n';

```

Listing 1.3: antlr.lexer

```

107 parser grammar ExprParser;
108 options { tokenVocab=ExprLexer; }
109
110 fisier
111     : intro bodypart outro;
112 bodypart
113     : break
114     | heading
115     | p
116     | bodypart break
117     | bodypart heading
118     | bodypart p
119     | ita
120     | bodypart ita
121     | bld
122     | bodypart bld;
123
124 break
125     : BR NL;
126

```

```

127 p
128     : STRING NL NL;
129
130 ita
131     : ITALIC STRING ITALIC NL;
132
133 bld : BOLD STRING BOLD NL;
134
135 heading
136     : HEADING1 STRING NL
137     | HEADING2 STRING NL
138     | HEADING3 STRING NL
139     | HEADING4 STRING NL
140     | HEADING5 STRING NL
141     | HEADING6 STRING NL;
142
143 intro;;
144
145 outro;;

```

Listing 1.4: antlr.parser

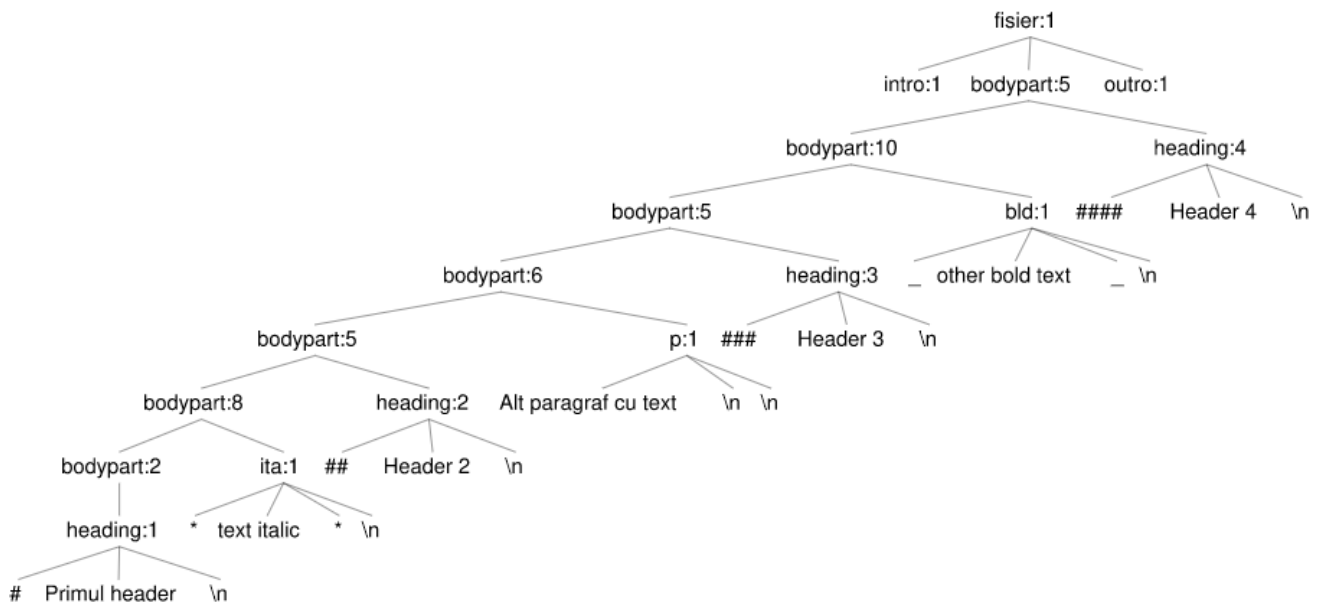
The screenshot displays the ANTLR online parser playground interface. On the left, the grammar from Listing 1.4 is loaded. The input field contains 'sample.expr'. The 'Run' button has been clicked, resulting in a parse tree visualization on the right.

**Parse Tree Hierarchy:**

```

graph TD
    fisier1[fisier:1] --> intro1[intro:1]
    fisier1 --> bodypart5_1[bodypart:5]
    fisier1 --> outro1[outro:1]
    intro1 --> Primul[Primul header]
    bodypart5_1 --> bodypart10[bodypart:10]
    bodypart5_1 --> heading4[heading:4]
    bodypart10 --> bodypart5_2[bodypart:5]
    bodypart10 --> bld1[bld:1]
    bodypart10 --> Header4[Header 4]
    bodypart5_2 --> bodypart6[bodypart:6]
    bodypart5_2 --> heading3[heading:3]
    bodypart6 --> bodypart5_3[bodypart:5]
    bodypart6 --> p1[p:1]
    bodypart5_3 --> bodypart8[bodypart:8]
    bodypart5_3 --> heading2[heading:2]
    bodypart8 --> bodypart2[bodypart:2]
    bodypart8 --> ita1[ita:1]
    bodypart8 --> Header2[Header 2]
    bodypart2 --> heading1[heading:1]
    bodypart2 --> textitalic1[text italic]
    ita1 --> textitalic2[text italic]
    heading1 --> textitalic3[text italic]
    
```





## 1.4 Rezultate

```

146 #Primul header
147 *text italic*
148 ## Header 2
149 Alt paragraf cu text
150
151 ### Header 3
152 _other bold text_
153 #### Header 4
154 **pepper**
155 asta e un paragraf
156
157 Un alt paragraf care contine text
158
159 ##### Header 5
160 Paragraf care contine descrierea documentului html

```

Listing 1.5: text.y

```

161 <!DOCTYPE html>
162 <html>
163 <body>
164   <h1>Primul header</h1>
165   <em>text italic</em>
166   <h2> Header 2 </h2>
167   <p>Alt paragraf cu text</p>
168   <h3> Header 3 </h3>
169   <strong>other bold text</strong>
170   <h4> Header 4 </h4>
171   <strong>pepper</strong>
172   <p>asta e un paragraf</p>
173   <p>Un alt paragraf care contine text</p>
174   <h5> Header 5</h5>
175   <p>Paragraf care contine descrierea documentului html</p>
176 </body>

```

## Listing 1.6: text.y

