Scraping The Simpsons Transcripts with R

Gaston Sanchez

Creative Commons Attribution Share-Alike 4.0 International CC BY-SA

The Simpsons Transcripts

The text of transcripts for The Simpsons are available at the Forever Dreaming Transcripts website:

https://transcripts.foreverdreaming.org/viewforum.php?f=431

https://transcripts.foreverdreaming.org/viewforum.php?f=431



URL of transcript episodes

Season 01, Episode 01: Simpsons Roasting on an Open Fire https://transcripts.foreverdreaming.org/viewtopic.php?f=431&t=21861

Season 01, Episode 02: Bart the Genius https://transcripts.foreverdreaming.org/viewtopic.php?f=431&t=21862

Season 01, Episode 03: Homer's Odyssey https://transcripts.foreverdreaming.org/viewtopic.php?f=431&t=21863

etc

Season 01, Episode 01

The first episode (from season 1) is:

"Simpsons Roasting on an Open Fire"

https://transcripts.foreverdreaming.org/viewtopic.php?f=431&t=21861





F.D. » Transcripts » S » The Simpsons

f y t 8+

Editor: SideshowBob

Print view

01x01 - Simpsons Roasting on an Open Fire

01x01 - Simpsons Roasting on an Open Fire

05/16/98 19:00

Marge: Ooh! Careful, Homer!

Homer: There's no time to be careful. We're late.

(tires screeching)

Class: J O little town of Bethlehem J

♪ How still we see thee lie ♪

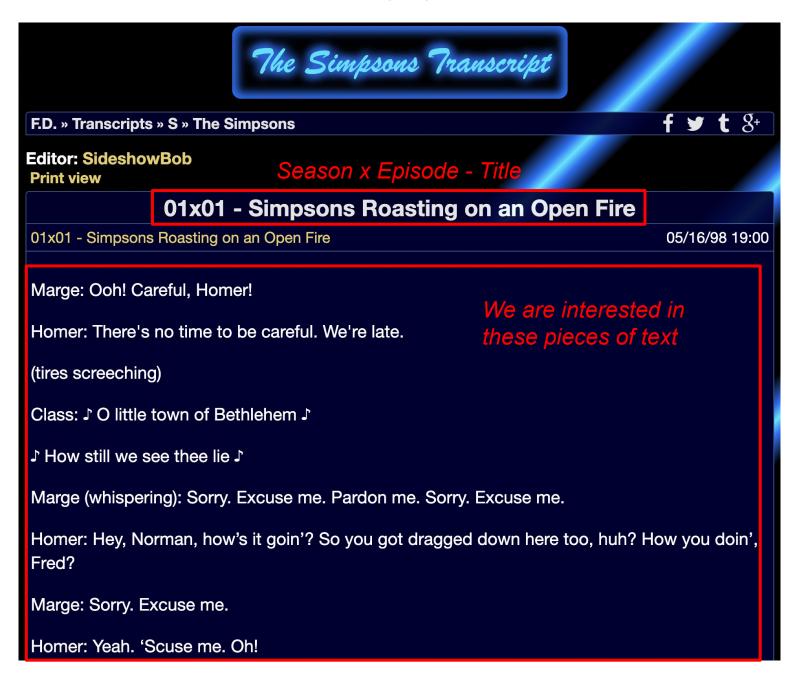
Marge (whispering): Sorry. Excuse me. Pardon me. Sorry. Excuse me.

Homer: Hey, Norman, how's it goin'? So you got dragged down here too, huh? How you doin',

Fred?

Marge: Sorry. Excuse me.

Homer: Yeah. 'Scuse me. Oh!



If we take a look at the source HTML code, we get:

```
Marge: Ooh! Careful, Homer!
121
422 Homer: There's no time to be careful. We're late.
(tires screeching)
  Class: ∫ 0 little town of Bethlehem ∫
  √p> I How still we see thee lie I
  Marge (whispering): Sorry. Excuse me. Pardon me. Sorry. Excuse me.
 Homer: Hey, Norman, how's it goin'? So you got dragged down here too, huh? How you
128 Marge: Sorry. Excuse me.
129 Homer: Yeah. 'Scuse me. Oh!
130 (yelps)
131 Homer: Pardon my galoshes. (laughs)
| 132 Class: J Are met in thee tonight J
133 (applause)
Skinner: (laughs) Wasn't that wonderful? And now, 'Santas of Many Lands,' as preser
135 Marge: Oh! Lisa's class.
137 (applause)
Boy: Merry Kurisumasu. I am Hotseiosha, a Japanese priest who acts like Santa Claus
139 (audience gasps)
40 Mr. Largo: Now presenting Lisa Simpson as Tawanga, the Santa Clause of the South Se
141 Homer: Ooh, it's Lisa! That's ours.
142 (drums b*at, natives chanting)
143 (applause)
Skinner: Ah, the fourth grade will now favor us with a melody- Uh, medley of holida
145 Class: ∫ Dashing through the snow ∫
146 ♪ In a one horse open sleigh ♪
147 ♪ O'er the fields we go ♪
|148|  1 Laughing all the way, ha ha ha |1 
149 ♪ Bells on bobtail ring ♪
|\langle p \rangle - \int \int (continues) \langle p \rangle
```

R Scripts & Files

R scripts

script1-scrape-episode-ids.R
script2-download-episode-html-files.R
script3-extract-transcript-lines.R
script4-assemble-output-table.R

```
script1-scrape-episode-ids.R
script2-download-episode-html-files.R
script3-extract-transcript-lines.R
script4-assemble-output-table.R
episode-ids.txt
simpsons-transcripts.txt
html_files/
   episode-21861.html
   episode-21862.html
   episode-73358.html
transcript_files/
   season-01-episode-01.txt
   season-01-episode-02.txt
   season-33-episode-22.txt
```

```
script1-scrape-episode-ids.R
script2-download-episode-html-files.R
script3-extract-transcript-lines.R
script4-assemble-output-table.R
episode-ids.txt ←---- output
simpsons-transcripts.txt
html_files/
   episode-21861.html
   episode-21862.html
   episode-73358.html
transcript_files/
   season-01-episode-01.txt
   season-01-episode-02.txt
   season-33-episode-22.txt
```

```
script1-scrape-episode-ids.R
      >script2-download-episode-html-files.R
       script3-extract-transcript-lines.R
input
       script4-assemble-output-table.R
                                                  output
      >episode-ids.txt
       simpsons-transcripts.txt
       html_files/
          episode-21861.html
          episode-21862.html
          episode-73358.html
       transcript_files/
          season-01-episode-01.txt
          season-01-episode-02.txt
          season-33-episode-22.txt
```

```
script1-scrape-episode-ids.R
       script2-download-episode-html-files.R
     >>script3-extract-transcript-lines.R
       script4-assemble-output-table.R
input ;
       episode-ids.txt
       simpsons-transcripts.txt
      `html_files/
          episode-21861.html
                                              output
          episode-21862.html
          episode-73358.html
       transcript_files/
          season-01-episode-01.txt
          season-01-episode-02.txt
          season-33-episode-22.txt
```

```
script1-scrape-episode-ids.R
     script2-download-episode-html-files.R
     script3-extract-transcript-lines.R
   *script4-assemble-output-table.R 
     episode-ids.txt
                                       output
     simpsons-transcripts.txt <----
     html_files/
        episode-21861.html
        episode-21862.html
input
        episode-73358.html
  ``-transcript_files/
        season-01-episode-01.txt
        season-01-episode-02.txt
        season-33-episode-22.txt
```

```
script1-scrape-episode-ids.R
script2-download-episode-html-files.R
script3-extract-transcript-lines.R
script4-assemble-output-table.R
episode-ids.txt
simpsons-transcripts.txt
html_files/
   episode-21861.html
   episode-21862.html
   episode-73358.html
transcript_files/
   season-01-episode-01.txt
   season-01-episode-02.txt
   season-33-episode-22.txt
```