

Программа htm2php

Преобразователь текста html в список строковых констант для php

Жуков И. Б.
ibzh@yandex.ru
<http://ibzh.16mb.com>
<http://ibzh.novhost.cf>

1. Программа *htm2php*. Программа *htm2php* преобразует html-текст в строки, пригодные для вставки в php. Исходно она предназначалась для обработки кода всевозможных счётчиков вроде liveinternet, yandex.метрика и т. п., а также других текстов, которые вам предлагают разместить на всех страницах вашего сайта. (Делать это, разумеется, проще всего с помощью php-скриптов.) Текст данной программы в интернете находится здесь.

2. Программа написана в стиле литературного программирования с использованием русифицированных средств. Чтобы её откомпилировать, вам понадобятся сами средства литературного программирования, которые можно найти по адресу

– <http://ibzh.16mb.com> или

– <http://ibzh.novhost.cf>,

а также компилятор языка паскаль, который можно скачать по адресу www.freepascal.com.

3. Программу будем писать по-русски, и для этого введём следующие тексты замены.

```

format без_конца ≡ until
define без_конца ≡ until (false)
format выйди_из_процедуры ≡ xclause
define выйди_из_процедуры ≡ exit
define выйди_из_функции ≡ exit
format выйди_из_функции ≡ xclause
format да ≡ true ;
define да ≡ true;
define делай ≡ do
format делай ≡ do
define если ≡ if
format если ≡ if
define заверши_программу ≡ halt
format заверши_программу ≡ xclause
define и_если ≡  $\wedge$ 
define иди_к ≡ goto
format иди_к ≡ goto
define или ≡  $\vee$ 
define иначе ≡ else
format иначе ≡ else
define конец ≡ end
format конец ≡ end
format константа ≡ const
define константа ≡ const
format константы ≡ const
define константы ≡ const
define метки ≡ label
format метки ≡ label
define начало ≡ begin
format начало ≡ begin
define не ≡  $\neg$ 
format нет ≡ false ;
define нет ≡ false;
format переменная ≡ var
define переменная ≡ var
define переменные ≡ var
format переменные ≡ var

```

4. Это продолжение определений предыдущего раздела.

```

define повторяй ≡ repeat
format повторяй ≡ repeat
define пока ≡ while
format пока ≡ while
define пока_не ≡ until
format пока_не ≡ until
format прерви_цикл ≡ xclause ;
define прерви_цикл ≡ break;
format программа ≡ program
define программа ≡ program
define процедура ≡ procedure
format процедура ≡ procedure
define то ≡ then
format то ≡ then
define функция ≡ function
format функция ≡ function

```

5. Программа должна из командной строки получить имя входного файла, обработать этот файл и записать итог в выходной файл. Имя выходного файла образуется из имени входного добавлением к нему расширения «.php».

Если программа отработала без ошибок, то операционной системе возвращается код завершения 0. Если была ошибка, то код завершения 1. (Код завершения можно проверить в командном файле, используя переменную окружения `%errorlevel%`.)

Ниже дан кусочек программы, выводящий справку по её использованию.

```

⟨ Выведи справку о программе и заверши её работу 5 ⟩ ≡
начало writeln; writeln(`Программа_преобразует_htm-файл_в_php-строки`); writeln;
writeln(`Использование:_htm2php_<имя_файла>`); writeln;
writeln(`Итог_записывается_в_файл_<имя_файла>.php`); writeln;
заверши_программу ;
конец;

```

Этот код используется в разделе 33.

6. Общий набросок программы выглядит следующим образом. Исходный файл считывается в буфер *a_in*, затем выполняется преобразование строк файла, итог которого помещается в буфер *a_out*. И, наконец этот буфер записывается на диск.

```

@{$COOPERATORS+@}
программа htm2php;
переменные ⟨ Общие переменные 10 ⟩
  ⟨ Процедуры и функции 8 ⟩
метки ⟨ Метки программы 28 ⟩
  начало ⟨ Открой входной файл 33 ⟩
  ⟨ Прочитай входной файл в буфер a_in 35 ⟩;
  ⟨ Закрой файл 36 ⟩
  ⟨ Выполни преобразование файла в строки. Итог сохрани в буфере в буфере a_out 26 ⟩
  ⟨ Открой выходной файл 37 ⟩
  ⟨ Запиши выходной файл 38 ⟩
  ⟨ Закрой выходной файл 39 ⟩
конец.

```

7. Между следующими командами будем размещать код для отладки программы. В обычной работе он не нужен, и поэтому он помещается в комментарии. Если потребуется что-то отлаживать, то нужно заменить определения **отладка** и **конец_отладки** на **begin** и **end**, соответственно.

```
define отладка ≡ @{ { begin }  
format отладка ≡ begin  
define конец_отладки ≡ @} { end; }  
format конец_отладки ≡ end ;
```

8. Когда программа обнаруживает ошибки, она вызывает процедуру *ошибка*, выводящую сообщение об ошибке и завершающую выполнение программы с кодом 1. Сообщение указывается строкой в аргументе функции.

```
⟨Процедуры и функции 8⟩ ≡  
процедура ошибка(s : string);  
  начало writeln; writeln(s);  
  заверши_программу (1);  
  конец;
```

См. также разделы 12 и 25.

Этот код используется в разделе 6.

9. Основная задача.

10. Сначала определим переменные, необходимые для обработки данных, а также процедуры доступа к данным и основных проверок. Указателями рассматриваемых мест в буферах *a_in* и *a_out* служат переменные *i* и *j*, соответственно. Наибольшие занятые места в буферах хранятся в переменных *s_in* и *s_out*. Наибольший размер входного буфера определяется константой *limit*. Выходной буфер делается немного больше, чтобы не выполнять проверки при вставке каждого знака в выходной буфер. (Дополнительные 10 мест в выходном буфере могут использоваться для размещения знаков, завершающих строку: кавычек, точки с запятой, точки, соединяющей константу со следующей строкой и др.)

Поскольку преобразование файла состоит в вставке в его содержимое дополнительных знаков, входной файл должен быть несколько меньше величины *limit*, иначе не хватит места в выходном буфере.

```
define limit = 10240 { наибольший размер входного файла }
```

⟨ Общие переменные 10 ⟩ ≡

```
a_in: array [0 .. limit] of byte; { здесь размещается входной файл }
```

```
a_out: array [0 .. limit + 10] of byte; { здесь составляется выходной файл }
```

```
s_in: dword = 0; { число входных символов }
```

```
s_out: dword = 0; { число выходных символов }
```

```
i, j: dword; { индексы в буферах a_in и a_out, соответственно }
```

См. также разделы 16 и 34.

Этот код используется в разделе 6.

11. Простейшие действия с буферами

```
define входной_знак ≡ a_in[i] { даёт текущий входной знак }
```

```
define есть_входные_знаки ≡ (i ≤ s_in) { проверка, что есть, что читать }
```

```
define нет_входных_знаков ≡ (i > s_in) { проверка обратного условия }
```

```
define есть_место_в_выходном_буфере ≡ (j < limit) { проверка, что есть, куда писать }
```

```
define нет_места_в_выходном_буфере ≡ (j ≥ limit) { проверка обратного условия }
```

```
define копируй_знак ≡
```

```
  начало a_out[j] ← входной_знак;
```

```
  отладка write (char (a_out[j]));
```

```
  конец_отладки i + = 1; j + = 1; { переносит знак из одного буфера в другой }
```

```
  конец
```

```
define помести_знак(#) ≡
```

```
  начало a_out[j] ← #;
```

```
  отладка write (char (a_out[j]));
```

```
  конец_отладки j + = 1; { Помещает знак в выходной буфер }
```

```
  конец
```

```
define следующий_входной_знак(#) ≡ a_in[#] { даёт k-ый входной знак }
```

```
define последний_входной_знак(#) ≡ (# = s_in) { проверяет, что k-ый знак был последним }
```

12. То, что мы подошли слишком близко к концу буфера, проверяет следующая процедура, она выдаёт сообщение об ошибке и завершает программу.

⟨ Процедуры и функции 8 ⟩ +≡

```
процедура проверь_место_в_выходном_буфере;
```

```
  начало если есть_место_в_выходном_буфере то выйди_из_процедуры ;
```

```
  ошибка ( `Выходной_буфер_слишком_мал` );
```

```
  конец;
```

13. Опишем суть задачи. На вход поступает HTML-текст, который должен построчно заключаться в кавычки. Текст помещается между кавычками один в один, кроме следующих символов, которые должны обрабатываться особым образом:

- одинарные кавычки;
- двойные кавычки;
- знаки обратной косой черты;
- знак "\$";
- переводы строки.

14. Строки в php могут заключаться или в одинарные кавычки, или в двойные. В последнем случае строки обрабатываются особым образом: в них раскрываются переменные, начинающиеся со знака "\$", и команды, начинающиеся со знака "\". Поэтому части текста HTML, содержащие указанные знаки, должны заключаться в одинарные кавычки.

Далее идёт подстановка, указывающая, что знак обрабатывается особым способом в двойных кавычках

```
define вх-знак-особый-в-двойных-кавычках ≡ ((входной-знак = "\")или(входной-знак = "$")или(входной-знак = ""))
```

15. С другой стороны, части строки, содержащие одинарные кавычки, должны заключаться в двойные. Нужно ли менять вид кавычек с одинарных на двойные, определяет следующая подстановка

```
define вх-знак-особый-в-обычных-кавычках ≡ (входной-знак = "'")
```

16. То, какие сейчас кавычки используются, будет определяться переменной *двойные-кавычки*

```
<Общие переменные 10> +=  
двойные-кавычки: boolean;
```

17. Если эта переменная истина, то HTML-текст сейчас помещается между двойными кавычками, иначе между одинарными. Переключение производится следующими подстановками

```
define теперь-двойные-кавычки ≡ двойные-кавычки ← true  
define теперь-обычные-кавычки ≡ двойные-кавычки ← false
```

18. Следующие подпрограммы меняют текущий вид кавычек.

```
define смени-кавычки-на-двойные ≡  
  begin помести-знак ("'); помести-знак ("."); помести-знак ("");  
  теперь-двойные-кавычки;  
  end  
define смени-кавычки-на-одинарные ≡  
  begin помести-знак (""); помести-знак ("."); помести-знак ("');  
  теперь-обычные-кавычки;  
  end
```

19. Рассмотрим теперь общий цикл обработки строки. Строка начинается со знака одинарной или двойной кавычки. Правило такое: если первый знак в исходном файле — одинарная кавычка, то преобразованная строка должна начинаться с двойной кавычки. Во всех прочих случаях — с одинарной. Этот кусочек кода обрабатывает начало непустой строки.

```
⟨ Обработай строку входного файла 19 ⟩ ≡
  если вх_знак_особый_в_обычных_кавычках то
    begin помести_знак("''"); теперь_двойные_кавычки;
    end
  иначе begin помести_знак("`"); теперь_обычные_кавычки;
  end;
```

См. также разделы 20, 21, 22, 23 и 24.

Этот код используется в разделе 26.

20. Следующий кусочек копирует данные до конца строки. Если встречается особый знак, то производится переключение на другой тип кавычки.

```
define конец_строки ≡ (входной_знак = 13)
```

```
⟨ Обработай строку входного файла 19 ⟩ +≡
  пока (неконец_строки) делай
    начало копируй_знак;
    если двойные_кавычки_и_если_вх_знак_особый_в_двойных_кавычках то
      смени_кавычки_на_одинарные
    иначе если (не_двойные_кавычки) и_если_вх_знак_особый_в_обычных_кавычках то
      смени_кавычки_на_двойные;
    если нет_входных_знаков то прерви_цикл ;
    проверь_место_в_выходном_буфере;
  конец;
```

21. В конце строки ставится завершающая кавычка соответствующего вида. Кроме того устанавливается переменная числа готовых для вывода символов *s_out*. (На самом деле *s_out* указывает на следующее место за последним готовым символом, так что в буфер *a_in[s_out]* нужно будет просто поместить точку с запятой перед выводом.)

```
⟨ Обработай строку входного файла 19 ⟩ +≡
  если двойные_кавычки то помести_знак("''")
  иначе помести_знак("`");
  s_out ← j;
  если нет_входных_знаков то ⟨ Иди к завершению обработки файла 29 ⟩
```

22. Сейчас мы указываем на знак с кодом 13, нам надо с него перейти на следующий. Также не забудем о том, что нужно пропустить символ с кодом 10, который во многих системах употребляется вместе со знаком 13.

```
⟨ Обработай строку входного файла 19 ⟩ +≡
  i + = 1;
  если (входной_знак = 10) то i + = 1; { передвинем указатель в a_in на одно место вправо }
```


23. Рассмотрим теперь, как должен обрабатываться конец строки. Если это не последняя строка файла, то нужно поставить точку и символы «`"\n"`». Последние символы заставят при обработке `phr` сохранить имевший место перевод на другую строку. Если следующая строка пустая, то логично поместить знаки `\n` в те же кавычки.

```

< Обработай строку входного файла 19 > +≡
  помести_знак (" "); помести_знак ("."); помести_знак (" "); помести_знак ("");
  повторяй помести_знак (" "); помести_знак ("n"); проверь_место_в_выходном_буфере;
    отладка writeln; write('i= ', i, ' s_in= ', s_in, ' s-in, ' s_out= ', j, ' s_out= ', s_out);
    writeln;
    конец_отладки;
  пока_не (не_пустая_строка);
  если нет_входных_знаков то <Иди к завершению обработки файла 29 >
  помести_знак (""); помести_знак (".");

```

24. Следующее, что нужно сделать, это поместить перевод на следующую строку в создаваемом `phr`-файле

```

< Обработай строку входного файла 19 > +≡
  помести_знак (13); помести_знак (10);

```

25. Теперь надо разобраться с тем, как определяется то, что следующая строка пустая. Делать мы это будем при помощи функции *пустая_строка*, которая будет просматривать знаки во входном буфере и возвращать *да*, если следующая строка пустая. Если строка из одних пробелов, то они все отбрасываются. Если дошли до конца строки, то перескакиваем с него на начало следующей строки, а если встретился символ, отличный от пробела и конца строки, то возвращаемся к началу строки.

```

define входной_знак_не_пробел ≡ ((входной_знак ≠ " ") и_если (входной_знак ≠ 9))

```

```

< Процедуры и функции 8 > +≡

```

```

функция пустая_строка: boolean;

```

```

  переменная k: word;

```

```

  начало k ← i; { Сохраним i }

```

```

  пустая_строка ← нет; { наиболее вероятный случай }

```

```

  пока (есть_входные_знаки) делай

```

```

    начало если конец_строки то

```

```

      начало i+ = 1;

```

```

      если входной_знак = 10 то i+ = 1; { перейдём к началу следующей строки }

```

```

      пустая_строка ← да;

```

```

      отладка writeln; writeln(' <пустая_строка: true; i= ', i, ' >');

```

```

      конец_отладки

```

```

      выйди_из_функции ;

```

```

      конец;

```

```

    если входной_знак_не_пробел то

```

```

      начало i ← k; { вернёмся к началу строки }

```

```

      отладка writeln; writeln(' <пустая_строка: false; i= ', i, ' >');

```

```

      конец_отладки

```

```

      выйди_из_функции ;

```

```

      конец;

```

```

    i+ = 1;

```

```

    конец;

```

```

  отладка writeln; writeln(' <пустая_строка: false; i= ', i, ' >');

```

```

  конец_отладки

```

```

конец;

```

26. Теперь можно описать следующий кусочек.

```

⟨Выполни преобразование файла в строки. Итог сохрани в буфере в буфере a_out 26⟩ ≡
  i ← 0; { Установим указатель на начало буфера }
  ⟨Пропусти пустые строки, если есть 27⟩
  повторяй ⟨Обработай строку входного файла 19⟩
  без_конца ;
завершение_обработки: ⟨Обработка конца вывода 30⟩

```

Этот код используется в разделе 6.

27. ⟨Пропусти пустые строки, если есть 27⟩ ≡
повторяй пока_не (*не_пустая_строка*);

Этот код используется в разделе 26.

28. Метки в программе — страшный сон для структурного объектно-ориентированного программиста.

```

⟨Метки программы 28⟩ ≡
  завершение_обработки;

```

Этот код используется в разделе 6.

29. ⟨Иди к завершению обработки файла 29⟩ ≡
goto *завершение_обработки*;

Этот код используется в разделах 21 и 23.

```

30. ⟨Обработка конца вывода 30⟩ ≡
  если s_out < 3 то ошибка('Входной_файл_не_содержит_данных');
  j ← s_out; помести_знак("");

```

Этот код используется в разделе 26.

31. На этом основная программа закончилась, остались только рутинные операции по открытию и закрытию файлов.

32. Работа с файлами.

33. Имя входного файла принимается из командной строки. В ней должно быть указано только одно имя

```

< Открой входной файл 33 > ≡
  если ParamCount < 1 то
    ошибка('Не_указан_входной_файл.' + #13#10 + 'Для_справки_наберите_htm2php-?');
  если ParamCount > 1 то ошибка('Слишком_много_параметров_во_входной_строке');
  если ParamStr(1) = '-?' то < Выведи справку о программе и заверши её работу 5 >
  assign(f, ParamStr(1)); @{$I-@}reset(f); @{$I+@}
  если IOResult ≠ 0 то ошибка('Не_могу_открыть_файл_"_ + ParamStr(1) +
    '_"'возможно,_он_имеет_недопустимое_имя');

```

Этот код используется в разделе 6.

34. Здесь нам понадобится переменная файлового типа. Будем рассматривать все файлы, как последовательности байт.

```

< Общие переменные 10 > +≡
f: file of byte;

```

35.

```

< Прочитай входной файл в буфер a_in 35 > ≡
  s_in ← 0;
  если eof(f) то ошибка('Входной_файл_пуст');
  пока (s_in < limit) делай
    начало read(f, a_in[s_in]);
    если eof(f) то прерви_цикл ;
    s_in + = 1;
  конец;
  если неeof(f) то ошибка('Входной_файл_слишком_длинный');

```

Этот код используется в разделе 6.

36. Следующий кусочек закрывает входной файл и, если выполняется отладка, выводит его содержимое на экран.

```

< Закрой файл 36 > ≡
  начало close(f);
  отладка writeln('Прочитан_файл:');
  for i ← 0 to s_in do write(char(a_in[i]));
  i ← 0; write('^'); writeln; writeln('Пишется_файл:');
  конец_отладки;
  конец;

```

Этот код используется в разделах 6 и 39.

37.

```

< Открой выходной файл 37 > ≡
  assign(f, ParamStr(1) + '.php'); @{$I-@}rewrite(f); @{$I+@}
  если IOResult ≠ 0 то ошибка('Не_могу_создать_файл_"_ + ParamStr(1) + '.php"');

```

Этот код используется в разделе 6.

38.

```
⟨Запиши выходной файл 38⟩ ≡  
  j ← 0;  
  пока (j ≤ s_out) делай  
    начало write(f, a_out[j]); j+ = 1;  
    конец;
```

Этот код используется в разделе 6.

39.

```
⟨Закрой выходной файл 39⟩ ≡  
  ⟨Закрой файл 36⟩  
  writeln; writeln(‘Готово. Итоговый_файл: <’ + ParamStr(1) + ‘.php>’);
```

Этот код используется в разделе 6.

40. Указатель.

без_конца: [3](#).*вх_знак_особый_в_двойных_кавычках:* [14](#), [20](#).*вх_знак_особый_в_обычных_кавычках:* [15](#), [19](#), [20](#).*входной_знак:* [11](#), [14](#), [15](#), [20](#), [22](#), [25](#).*входной_знак_не_пробел:* [25](#).**выйди_из_процедуры:** [3](#).**выйди_из_функции:** [3](#).*да:* [3](#), [25](#).*двойные_кавычки:* [16](#), [17](#), [20](#), [21](#).**делай:** [3](#).**если:** [3](#).*есть_входные_знаки:* [11](#), [25](#).*есть_место_в_выходном_буфере:* [11](#), [12](#).*завершение_обработки:* [26](#), [28](#), [29](#).**заверши_программу:** [3](#).*и_если:* [3](#), [20](#), [25](#).**иди_к:** [3](#).*или:* [3](#), [14](#).**иначе:** [3](#).**конец:** [3](#).**конец_отладки:** [7](#).*конец_строки:* [20](#), [25](#).**константа:** [3](#).**константы:** [3](#).*копируй_знак:* [11](#), [20](#).**метки:** [3](#).**начало:** [3](#).*не:* [3](#), [20](#), [23](#), [27](#), [35](#).*нет:* [3](#), [25](#).*нет_входных_знаков:* [11](#), [20](#), [21](#), [23](#).*нет_места_в_выходном_буфере:* [11](#).**отладка:** [7](#).*ошибка:* [8](#), [12](#), [30](#), [33](#), [35](#), [37](#).**переменная:** [3](#).**переменные:** [3](#).**повторяй:** [4](#).**пока:** [4](#).**пока_не:** [4](#).*помести_знак:* [11](#), [18](#), [19](#), [21](#), [23](#), [24](#), [30](#).*последний_входной_знак:* [11](#).**прерви_цикл:** [4](#).*проверь_место_в_выходном_буфере:* [12](#), [20](#), [23](#).**программа:** [4](#).**процедура:** [4](#).*пустая_строка:* [23](#), [25](#), [27](#).*следующий_входной_знак:* [11](#).*смени_кавычки_на_двойные:* [18](#), [20](#).*смени_кавычки_на_одинарные:* [18](#), [20](#).*теперь_двойные_кавычки:* [17](#), [18](#), [19](#).*теперь_обычные_кавычки:* [17](#), [18](#), [19](#).**то:** [4](#).**функция:** [4](#).*a_in:* [6](#), [10](#), [11](#), [21](#), [22](#), [35](#), [36](#).*a_out:* [6](#), [10](#), [11](#), [38](#).*assign:* [33](#), [37](#).**begin:** [3](#), [7](#).*boolean:* [16](#), [25](#).*break:* [4](#).*byte:* [10](#), [34](#).*char:* [11](#), [36](#).*close:* [36](#).**const:** [3](#).**COPERATORS:** [6](#).**do:** [3](#).*dword:* [10](#).**else:** [3](#).**end:** [3](#), [7](#).*eof:* [35](#).*exit:* [3](#).*f:* [34](#).*false:* [3](#), [17](#).**function:** [4](#).**goto:** [3](#).*halt:* [3](#).*htm2php:* [1](#), [6](#).*i:* [10](#).**if:** [3](#).*IOResult:* [33](#), [37](#).*k:* [25](#).**label:** [3](#).*limit:* [10](#), [11](#), [35](#).*ParamCount:* [33](#).*ParamStr:* [33](#), [37](#), [39](#).**procedure:** [4](#).**program:** [4](#).*read:* [35](#).**repeat:** [4](#).*reset:* [33](#).*rewrite:* [37](#).*s_in:* [10](#), [11](#), [23](#), [35](#), [36](#).*s_out:* [10](#), [21](#), [23](#), [30](#), [38](#).*string:* [8](#).**then:** [4](#).*true:* [3](#), [17](#).**until:** [3](#), [4](#).**var:** [3](#).**while:** [4](#).*word:* [25](#).*write:* [11](#), [23](#), [36](#), [38](#).*writeln:* [5](#), [8](#), [23](#), [25](#), [36](#), [39](#).**xclause:** [3](#), [4](#).

- ⟨ Выведи справку о программе и заверши её работу 5 ⟩ Используется в разделе 33.
- ⟨ Выполни преобразование файла в строки. Итог сохрани в буфере в буфере *a_out* 26 ⟩
Используется в разделе 6.
- ⟨ Закрой выходной файл 39 ⟩ Используется в разделе 6.
- ⟨ Закрой файл 36 ⟩ Используется в разделах 6 и 39.
- ⟨ Запиши выходной файл 38 ⟩ Используется в разделе 6.
- ⟨ Иди к завершению обработки файла 29 ⟩ Используется в разделах 21 и 23.
- ⟨ Метки программы 28 ⟩ Используется в разделе 6.
- ⟨ Обработай строку входного файла 19, 20, 21, 22, 23, 24 ⟩ Используется в разделе 26.
- ⟨ Обработка конца вывода 30 ⟩ Используется в разделе 26.
- ⟨ Общие переменные 10, 16, 34 ⟩ Используется в разделе 6.
- ⟨ Открой входной файл 33 ⟩ Используется в разделе 6.
- ⟨ Открой выходной файл 37 ⟩ Используется в разделе 6.
- ⟨ Пропусти пустые строки, если есть 27 ⟩ Используется в разделе 26.
- ⟨ Процедуры и функции 8, 12, 25 ⟩ Используется в разделе 6.
- ⟨ Прочитай входной файл в буфер *a_in* 35 ⟩ Используется в разделе 6.

	Раздел	Стр.
Программа <i>htm2php</i>	1	3
Основная задача	9	6
Работа с файлами	32	11
Указатель	40	13