

2017年5月10日 由YUER

为htmlpurifier订制xss过滤

xss简介

我们每天读的公众号文章，其内容是小编通过富文本编辑器生成的。内容发布者在富文本编辑器里可视化的编写文字与图片并进行简单的排版，背后其实是编辑器动态的生成各种html标签，最终文章对应的完整html会被提交到服务端保存。当读者打开文章时，服务端将保存的html直接返回给浏览器进行渲染。

读者直接获取并展现小编提交的html，这是极其危险的，万一html里有一段<script>包裹的Js代码，那读者的信息就可以被小编通过Js代码盗走，这种攻击手段叫做xss。当然，像<script>这种xss手法是很好防御的，只需要检查小编提交的html过滤掉<script>标签即可，但是xss远没有这么简单。恶意的小编会利用错乱的<、>、引号等导致最终读者的页面出现意想不到的标签闭合效果，或者利用onclick，onerror等标签的事件函数来实现js代码的执行，等等诸如此类的注入手法。

那么xss过滤就是用来检查小编提交的html，将存在安全隐患的代码移除掉，将不规范的代码整理或移除，最终留下一个安全可靠的html，安全到可以直接展现给读者。

htmlpurifier基础使用

既然xss的注入手段各式各样，光是事件函数的类型都数不过来，因此不如找一个历史积累丰富的xss过滤开源库，那就是著名的htmlpurifier了，像有名的Yi2框架也是内置了这个库。

如果你是composer项目，可以通过composer安装与自动加载。如果是比较老的项目，可以下载standalone版本，并通过require_once引入。

基础用法不需要任何配置，直接将小编提交的html传给htmlpurifier即可完成过滤：

```
1 $config = HTMLPurifier_Config::createDefault();
2 $purifier = new HTMLPurifier($config);
3 $clean_html = $purifier->purify($dirty_html);
```

htmlpurifier添加标签和属性

这个库遵循HTML标准，会保留常见的各种HTML标签，而删除库中未定义的标签。对于标签属性也是如此，只会保留合法的HTML定义的标准属性。通常来说，你不必对它做什么订制就可以正常工作了，但是也不排除特殊需求。

一些情况下，你会希望保留某些自定义属性，或者保留某些自定义的标签，这时候必须订制htmlpurifier。订制并不是很复杂，可以参考官方文档来学习如何订制标签和属性，网上鲜有博客说明这一块的具体做法。

我遇到的情况是，公司的某个富文本中包含<dir>标签被htmlpurifier删除了，我查了一下HTML标准获知，这个标签是符合标准的，但是不推荐使用了。没办法，我必须订制它，避免htmlpurifier将其过滤：

```
1 // 默认配置
2 $config = HTMLPurifier_Config::createDefault();
3 // 自定义配置 http://htmlpurifier.org/docs/enduser-customize.html
4 $config->set('HTML.DefinitionID', 'smzdm library version');
5 $config->set('HTML.DefinitionRev', 2);
6 if ($def = $config->maybeGetRawHTMLDefinition()) {
7     // 允许卡片的<dir>标签
8     $def->addElement(
9         'dir',
10         'Block',
11         'Flow',
12         'Common',
13         [
14             'res-data-id' => 'Number'
15         ]
16     );
17 }
18
19 $this->html_purifier = new HTMLPurifier($config);
```

这里多了2个set操作，第一个是设置这个配置的名称，第二个是配置的版本，它俩会作为配置的缓存key，最终这份配置会被缓存到磁盘到一个文件里。

maybeGetRawHTMLDefinition()会检查名称+版本的缓存文件是否存在，如果已经存在则返回的\$def为null并直接使用缓存文件的配置，否则会执行你的配置并缓存到文件中去。因此，如果你正在调试各种配置项，那么每次修改完配置代码都应该修改版本号来避免加载缓存的配置，以便看到实时的修改结果。

另外一种调试期间防止缓存的方法，是将maybeGetRawHTMLDefinition替换成永远不加载缓存的API：

```
1 // 默认配置
2 $config = HTMLPurifier_Config::createDefault();
3 // 自定义配置 http://htmlpurifier.org/docs/enduser-customize.html
4 $config->set('HTML.DefinitionID', 'smzdm library version');
5 $config->set('HTML.DefinitionRev', 1);
6 $config->set('Cache.DefinitionImpl', null); // 清理配置缓存,上线时关掉这句
7 $def = $config->getHTMLDefinition(true);
8 // 允许卡片的<dir>标签
9 $def->addElement(
10     'dir',
11     'Block',
12     'Flow',
13     'Common',
14     [
15         'res-data-id' => 'Number'
16     ]
17 );
18
19 $this->html_purifier = new HTMLPurifier($config);
```

这可以方便你调试配置时立即看到效果，但是上线时请用第一份代码，因为它会生成并且只生成一次缓存，而第二份代码永远不会生成缓存，性能差并且会发出代码警告。

我又通过addElement订制了一个<dir>标签。

它是一个块级元素（Block），我们知道合理的HTML规范是块级元素内可以有块级元素与内联元素，而内联元素内只能有内联元素。

Flow指定了<dir>允许孩子是块级元素或者内联元素，这和HTML规范一致。

Common指定了<dir>可以包含一些常见的属性：id, style, class, title and lang..

最后一个参数是为<dir>添加了自定义属性，这样res-data-id属性就不会被htmlpurifier过滤掉了，它的合法值是一个整形，否则将被过滤掉。

封装htmlpurifier

公司的富文本框做了一些订制，比如允许嵌入embed视频。为了方便配置和管理Htmlpurifier，我仿照这个开源项目对其进行了进一步封装如下。注意，该开源项目在配置文件缓存使用方面存在BUG，我已对其修复。

```
1 <?php
2 /**
3  * require_once(__DIR__ . '/HTMLPurifier.standalone.php');
4  *
5  * /**
6  *  * Class purifier
7  *  *
8  *  * 富文本过滤器
9  *  *
10 *  * 封装 http://htmlpurifier.org/docs
11 *  */
12
13 class Purifier
14 {
15     /**
16      * @var config
17      */
18     protected $config;
19
20     /**
21      * @var HTMLPurifier
22      */
23     protected $purifier;
24
25     /**
26      * Constructor
27      *
28      * @throws Exception
29      */
30     public function __construct()
31     {
32         $this->config = require(__DIR__ . '/config/purifier.php');
33         $this->setUp();
34     }
35
36     /**
37      * @param $key
38      * @param null $default
39      * @return null
40      */
41     private function config($key, $default = null) {
42         $paths = explode('.', $key);
43
44         $cur_config = $this->config;
45         foreach ($paths as $path) {
46             if (!isset($cur_config[$path])) {
47                 return $default;
48             }
49             $cur_config = $cur_config[$path];
50         }
51         return $cur_config;
52     }
53
54     /**
55      * Setup
56      *
57      * @throws Exception
58      */
59     private function setUp()
60     {
61         // Create a new configuration object
62         $config = HTMLPurifier_Config::createDefault();
63
64         // 基础配置
65         $config->loadArray($this->getConfig());
66
67         // 订制配置
68         $definition = $this->config('settings.custom_definition');
69         if (empty($definition)) {
70             $config->set('HTML.DefinitionID', $definition['id']);
71             $config->set('HTML.DefinitionRev', $definition['rev']);
72             // Enable debug mode
73             if (isset($definition['debug']) || $definition['debug']) {
74                 $config->set('Cache.DefinitionImpl', null);
75             }
76             // 优先加载缓存的配置，如果不存在就执行订制代码
77             if ($def = $config->maybeGetRawHTMLDefinition()) {
78                 $this->addCustomDefinition($definition, $def);
79             }
80         }
81
82         // Create HTMLPurifier object
83         $this->purifier = new HTMLPurifier($this->configure($config));
84     }
85
86     /**
87      * Add a custom definition
88      *
89      * @see http://htmlpurifier.org/docs/enduser-customize.html
90      * @param array $definitionConfig
91      * @param HTML_Purifier_Config $configObject Defaults to using default config
92      * @return HTML_Purifier_Config $configObject
93      */
94     private function addCustomDefinition(array $definitionConfig, $defObj = null)
95     {
96         // Create the definition attributes
97         if (empty($definitionConfig['attributes'])) {
98             $this->addCustomAttributes($definitionConfig['attributes'], $defObj);
99         }
100
101         // Create the definition elements
102         if (empty($definitionConfig['elements'])) {
103             $this->addCustomElements($definitionConfig['elements'], $defObj);
104         }
105     }
106
107     /**
108      * Add provided attributes to the provided definition
109      *
110      * @param array $attributes
111      * @param HTMLPurifier_HTMLDefinition $definition
112      * @return HTMLPurifier_HTMLDefinition $definition
113      */
114     private function addCustomAttributes(array $attributes, $definition)
115     {
116         foreach ($attributes as $attribute) {
117             // Get configuration of attribute
118             $required = !empty($attribute[3]) ? true : false;
119             $onElement = $attribute[0];
120             $attrName = $attribute[1];
121             $validValues = $attribute[2];
122
123             $definition->addAttribute($onElement, $attrName, $validValues);
124         }
125
126         return $definition;
127     }
128
129     /**
130      * Add provided elements to the provided definition
131      *
132      * @param array $elements
133      * @param HTMLPurifier_HTMLDefinition $definition
134      * @return HTMLPurifier_HTMLDefinition $definition
135      */
136     private function addCustomElements(array $elements, $definition)
137     {
138         foreach ($elements as $element) {
139             // Get configuration of element
140             $name = $element[0];
141             $contentSet = $element[1];
142             $allowedChildren = $element[2];
143             $attributeCollection = $element[3];
144             $attributes = !empty($element[4]) ? $element[4] : null;
145
146             if (empty($attributes)) {
147                 $definition->addElement($name, $contentSet, $allowedChildren, $attributeCollection, $attributes);
148             } else {
149                 $definition->addElement($name, $contentSet, $allowedChildren, $attributeCollection);
150             }
151         }
152     }
153
154     /**
155      * @param HTMLPurifier_Config $config
156      * @return HTMLPurifier_Config
157      */
158     protected function configure(HTMLPurifier_Config $config)
159     {
160         return $config;
161     }
162
163     /**
164      * @return array|null
165      */
166     protected function getConfig()
167     {
168         $default_config = [];
169         $default_config['Core.Encoding'] = $this->config('encoding');
170
171         $config = $this->config('settings.default');
172
173         if (is_array($config)) {
174             $config = [];
175         }
176
177         $config = $default_config + $config;
178
179         return $config;
180     }
181
182     /**
183      * @param $dirty
184      * @return mixed
185      */
186     public function purify($dirty)
187     {
188         return $this->purifier->purify($dirty);
189     }
190 }
```

配置文件可以这样来写：

```
1 <?php
2 /**
3  * Ok, glad you are here
4  *
5  * first we get a config instance, and set the settings
6  *
7  * $config = HTMLPurifier_Config::createDefault();
8  * $config->set('Core.Encoding', $this->config->get('purifier.encoding'));
9  * $config->set('Core.SerializerPath', $this->config->get('purifier.serializePath'));
10 * if ( ! $this->config->get('purifier.finalize') ) {
11 *     $config->autoFinalize = false;
12 * }
13 * $config->loadArray($this->getConfig());
14 *
15 * You must NOT delete the default settings
16 * anything in settings should be compacted with params that needed to instance HTMLPurifier_Config.
17 *
18 * @link http://htmlpurifier.org/live/configdoc/plain.html
19 */
20
21 return [
22     'encoding' => 'UTF-8',
23     'settings' => [
24         // 默认配置
25         'default' => [
26             // 采用HTMLPurifier默认即可，一旦配置将完全覆盖HTMLPurifier的默认值，自定义必须覆盖所有tag和attr
27
28             // 'HTML.Doctype' => 'HTML 4.01 Transitional',
29             // 'HTML.Allowed' => 'div,b,strong,i,em,u,a[href=title],ul,ol,li,p[style],br,span[style],img[width,height,alt,src]',
30             // 'CSS.AllowedProperties' => 'font,font-size,font-weight,font-style,font-family,text-decoration,padding-left,color,background-color,text-align',
31             // 给文字包裹<p>标签
32             // 'AutoFormat.AutoParagraph' => true,
33             // 删除没有孩子的标签
34             // 'AutoFormat.RemoveEmpty' => true,
35
36             'HTML.SafeEmbed' => true, // 允许嵌入视频
37             'HTML.SafeIframe' => true, // 允许iframe
38             'URI.SafeIframeRegexp' => '%http://(.+?.youku.com/|.+?.tudou.com/|.+?.56.com%)/%', // 允许iframe的src校验
39             'Attr.DefaultImageAlt' => '',
40
41             'custom_definition' => [
42                 'id' => 'my-definitions',
43                 'rev' => 2,
44                 'debug' => false,
45                 'elements' => [
46                     // http://developers.whatwg.org/sections.html
47                     ['section', 'Block', 'Flow', 'Common'],
48                     ['nav', 'Block', 'Flow', 'Common'],
49                     ['article', 'Block', 'Flow', 'Common'],
50                     ['aside', 'Block', 'Flow', 'Common'],
51                     ['header', 'Block', 'Flow', 'Common'],
52                     ['footer', 'Block', 'Flow', 'Common'],
53
54                     // Content model actually excludes several tags, not modelled here
55                     ['address', 'Block', 'Flow', 'Common'],
56                     ['hgroup', 'Block', 'Required: h1 | h2 | h3 | h4 | h5 | h6', 'Common'],
57
58                     // http://developers.whatwg.org/grouping-content.html
59                     ['figure', 'Block', 'Optional: (Content, Flow) | (Flow, figcaption) | Flow', 'Common'],
60                     ['figcaption', 'Inline', 'Flow', 'Common'],
61
62                     // http://developers.whatwg.org/the-video-element.html#the-video-element
63                     ['video', 'Block', 'Optional: (source, Flow) | (Flow, source) | Flow', 'Common'],
64                     ['src' => 'URI',
65
66                     'type' => 'Text',
67                     'width' => 'Length',
68                     'height' => 'Length',
69                     'poster' => 'URI',
70                     'preload' => 'Enum#auto,metadata,none',
71                     'controls' => 'Bool',
72                     ],
73                     ['source', 'Block', 'Flow', 'Common', [
74                         'res-data-id' => 'Number',
75                     ]],
76                     'attributes' => [
77                         ['iframe', 'allowfullscreen', 'Bool'],
78                         ['table', 'height', 'Text'],
79                         ['td', 'border', 'Text'],
80                         ['th', 'border', 'Text'],
81                         ['tr', 'width', 'Text'],
82                         ['tr', 'height', 'Text'],
83                         ['tr', 'border', 'Text'],
84                     ],
85                 ],
86             ],
87             // http://developers.whatwg.org/text-level-semantic.html
88             ['s', 'Inline', 'Inline', 'Common'],
89             ['var', 'Inline', 'Inline', 'Common'],
90             ['sub', 'Inline', 'Inline', 'Common'],
91             ['sup', 'Inline', 'Inline', 'Common'],
92             ['mark', 'Inline', 'Inline', 'Common'],
93             ['wbr', 'Inline', 'Empty', 'Core'],
94
95             // http://developers.whatwg.org/edits.html
96             ['ins', 'Block', 'Flow', 'Common', ['cite' => 'URI', 'datetime' => 'CDATA']],
97             ['del', 'Block', 'Flow', 'Common', ['cite' => 'URI', 'datetime' => 'CDATA']],
98
99             // 允许使用dir
100             ['dir', 'Block', 'Flow', 'Common', [
101                 'res-data-id' => 'Number',
102             ]],
103         ],
104     ],
105 ];
```

这份代码可以正常工作，通过配置即可改变htmlpurifier的行为，便于你探索与调试。

本文抛砖引玉，如果有疑惑就认真读一下官方文档，多测试多尝试。