

OcrKing 配合火车采集器简易说明

应管理 **真视觉** 之邀请,特做此简单说明文档,时间比较仓促,不详之处请多多见谅!首先请在本地或远程都搭建好 php 环境(推荐 xampp 或 wamp),安装即可用,并下载 OcrKing.php 脚本,请用文本编译器 notepad++或系统记事本打开 OcrKing.php,文件内会有以下内容,请获取 key 后替换下面的 **你的 apiKey** 并保存脚本

```
//本识别服务为免费,没有 apiKey 可以 key 为标题
// 任意内容为正文发邮件到 ok@ocrking.com 获取
//可能会有延迟,请勿重复发送
//授权 apiKey,请注意区分大小写
define ('API_KEY','你的 apiKey');
```

将保存好的脚本上传到你的 php 空间,然后访问

<http://domain.com/OcrKing.php?u=http://www.e-fa.cn/extend/image.php?auth=EXcQdVFjIkCOPxd1DW4wAF5iNAUJPQ>

如果返回结果为 **15201124022** 就说明成功了

配合火车,适合各个版本的火车采集器,主要利用了**多页功能**进行识别

请下载 **OcrKing** 演示(v8 规则,其它版本类似) 规则导入,查看多页设置仿照设置即可,如要采集的内容如下图

```
<ul class="contentUl">
<div class="lianxi_1">手机! </div>
<div class="lianxi_1">电话! </div>
<div class="lianxi_1">所在地: 北京</div>
<div class="lianxi_1">地址: 北京市海淀区上庄工业园区A32号</div>
<div class="lianxi_1"><a href="http://www.e-fa.cn/com/dianchi0530/" target="_blank" title="北京通隆恒盛科技有限公司">http://www.e-fa.cn/com/d
</ul>
<div class="clear"></div></div>
```

手机电话图片代码为

手机:

电话:

标签编辑

标签名: ☐ 该标签循环匹配 ☐ 该标签在分页中匹配 ☐ 从网址中采集

☒ 通过采集得到数据 ☐ 自定义固定格式的数据

提取数据方式
☒ 前后截取 ☐ 正则提取 ☐ 可视化提取 ☐ 正文提取 ☐ 标签组合 所属多页:

开始字符串: 结束字符串:

数据处理

添加 删除 清空

文件下载选项

☐ 将相对地址补全为绝对地址 ☐ 下载图片
☐ 探测文件真实地址但不下载 ☐ 探测文件并下载

文件地址必含:
 文件保存目录:
 文件保存格式:

确定 取消

第一步: 采集网址规则 | 第二步: 采集内容规则 | 第三步: 发布内容设置 | 文件保存及部分高级设置

页面内容标签定义 (规则普通编辑模式)

操作	标签名	开始...	结束...
添加	标题	<title>	</title>
修改	手机	<ocr...	</ocr...
复制	电话	<ocr...	</ocr...
粘贴			
删除			
导入			
多页管理			

切换到无限级多页规则编辑模式

分页获取规则 | 标签循环处理 | 其他设置

每一个标签循环匹配项

☐ 添加为新记录 ☒ 用分隔符连接在上条记录后

规则测试

典型页面: 测

请求页面 默认页 http://www.e-fa.cn/sell/201410/15/info_8481861.html
 请求页面 手机 http://127.0.0.1:90/0crKing.php?u=http://www.e-fa.cn/extend/image.php?auth=EXcQdVFjIkcoPxdIDW4wAF5iNAUJPQ
 请求页面 电话 http://127.0.0.1:90/0crKing.php?u=http://www.e-fa.cn/extend/image.php?auth=LkksTXtLE2syByZDUTATJH1EGSNDdkQo

【标题】: 赛特蓄电池MSE-400_中国易发网
 【手机】: 15201124022
 【电话】: 01056035938

手机: 15201124022
 电话: 010-56035938
 所在地: 北京