# website-comapre

stardiviner

2018 年 11 月 23 日

## 目录

# 1 INPROGRESS website-compare [37/51]

**DEADLINE:** *<2018-09-26 Wed>*  **SCHEDULED:** *<2018-09-22 Sat>*

website-compare/project.clj

# 2 新旧网站的基本信息

- 新网站 `http://www.sxszjzx.com`, `http://zjzx.sxsedu.net`

- 旧网站 `http://www.sxti.zj.cn`

内网解析 `http://www.sxszjzx.com/` 是解析到旧网站的。外网解析 `http://www.sxszjzx.com/` 是解析到新网站的。如果你解析到外网，你需要修改 DNS 到 172.18.0.5

## 2.1 TODO 外网无法访问旧网站

**DEADLINE:** *<2018-11-19 Mon>*

`http://bbs.360.cn/thread-6776749-1-1.html`

Modify the host file:

```
1 sudo echo "192.168.1.1 www.sxti.zj.cn" >> /etc/hosts
```

表 1: Clock summary at *[2018-11-22 Thu 21:30]*

| Headline | Time | |
|---|---|---|
| **Total time** | **4d 8:53** | |
| 写一个简单的爬虫加交叉对比的脚本来比较新旧网站的内容差别 [37/51] | 4d 8:53 | |
| 代码架构 | | 0:29 |
| website 新旧网站使用相同 API，不同的实现 [2/2] | | 3:55 |
| website-old.clj HTML elements | | |
| website-new.clj HTML elements | | |
| crawler.clj 获取所有文章的方法 [15/17] | | 1d 2:11 |
| benchmarking the HTTP URL request... | | |
| 按照网站的导航栏目录入口获取相应类目下的所有文章 | | 1 |
| old | | |
| get navigator bar links | | |
| parse sub-nav page to get all listed... | | |
| iterate all result pages | | |
| new | | |
| get navigator bar links | | |
| parse sub-nav page to get all listed... | | |
| iterate all result pages | | |
| benchmark new website two entry URL | | |
| crawler.clj 爬虫存储提取连接的数据库 [5/5] | | 1:57 |
| Redis | | |
| extract all links in nav bar | | |
| save all links to Redis | | |
| crawler.clj 爬取页面的内容 [3/6] | | 0:49 |
| scrape the article | | |
| old | | |
| new | | |
| process the crawled article content... | | |
| strip some special characters and spaces | | |
| store.clj 爬虫内容存储到数据库 [3/4] | | 1:03 |
| SQLite | | |
| save title and content to SQLite | | |
| 比较上传视频 | | 0:50 |
| 检查视频是否可以播放 [4/5] | | |

2

# 3 DONE 代码架构

**CLOSED:** *[2018-11-16 Fri 15:41]*

# 4 DONE website 新旧网站使用相同 API，不同的实现 [2/2]

**CLOSED:** *[2018-11-19 Mon 19:11]*
Use different namespace, same API function names.

## 4.1 DONE website-old.clj HTML elements

**CLOSED:** *[2018-11-19 Mon 19:05]*

### 4.1.1 nav

```
1 (def selector-nav [:div.head_3 :ul#head_nav])
```

```
1 <div class="head_3">
2   <ul id="head_nav">
3
4   <li class="nav_active"><a href="/"><span>首页</span></a></li>
5   <li class=""><a href="javascript:void(0)"><span>学校概况</span></a>
6     <ul>
7       <li><a href="/html/school/about.html">学校简介</a></li>
8       <li><a href="/e/action/ListInfo/?classid=13">现任领导</a></li>
9       <li><a href="/e/action/ListInfo/?classid=283">名优教师</a></li>
10      <li><a href="/html/school/office.html">处室联系</a></li>
```

```
11        <li><a href="/e/action/ListInfo/?classid=14">校园风景</a></li>

12        <li><a href="/e/action/ListInfo/?classid=35">楼层分布</a></li>

13        <li><a href="/html/school/map.html">交通地图</a></li>

14        <li><a href="/e/action/ListInfo/?classid=15">学校荣誉</a></li>

15        <li><a href="/e/action/ListInfo/?classid=377">媒体关注</a></li>

16      </ul>

17    </li>

18    <li class=""><a href="javascript:void(0)"><span>新闻频道</span></a>

19      <ul>

20        <li><a href="/e/action/ListInfo/?classid=33">校园新闻</a></li>

21        <li><a href="/e/action/ListInfo/?classid=34">系部新闻</a></li>

22        <li><a href="/e/action/ListInfo/?classid=7">对外交流</a></li>

23        <li><a href="/e/action/ListInfo/?classid=8">学校荣誉</a></li>

24        <li><a href="/e/action/ListInfo/?classid=9">师生荣誉</a></li>

25        <li><a href="/e/action/ListInfo/?classid=10">校园视频</a></li>

26        <li><a href="/e/action/ListInfo/?classid=11">二周安排</a></li>

27        <li><a href="/e/action/ListInfo/?classid=12">每周寄语</a></li>

28      </ul>

29    </li>

30    <li><a href="javascript:void(0)"><span>下载频道</span></a>

31      <ul>

32        <li><a href="/e/action/ListInfo/?classid=16">党政办</a></li>

33        <li><a href="/e/action/ListInfo/?classid=17">教务处</a></li>

34        <li><a href="/e/action/ListInfo/?classid=18">德育团委</a></li>

35        <li><a href="/e/action/ListInfo/?classid=53">总务处</a></li>

36        <li><a href="/e/action/ListInfo/?classid=54">实习处</a></li>

37        <li><a href="/e/action/ListInfo/?classid=56">安保处</a></li>

38        <li><a href="/e/action/ListInfo/?classid=175">教科室</a></li>

39        <li><a href="/e/action/ListInfo/?classid=57">教学软件</a></li>

40        <li><a href="/e/action/ListInfo/?classid=58">其它</a></li>

41      </ul>
```

```html
42    </li>
43    <li><a href="javascript:void(0)"><span>处室网站</span></a>
44      <ul>
45        <li><a href="/e/action/ListInfo/?classid=37">党政办</a></li>
46        <li><a href="/e/action/ListInfo/?classid=38">教务处</a></li>
47        <li><a href="/e/action/ListInfo/?classid=39">德育团委</a></li>
48        <li><a href="/e/action/ListInfo/?classid=40">总务处</a></li>
49        <li><a href="/e/action/ListInfo/?classid=41">实习处</a></li>
50        <li><a href="/e/action/ListInfo/?classid=43">安保处</a></li>
51        <li><a href="/e/action/ListInfo/?classid=104">教科室</a></li>
52        <li><a href="/e/action/ListInfo/?classid=44">服务公司</a></li>
53      </ul>
54    </li>
55    <li><a href="javascript:void(0)"><span>系部网站</span></a>
56      <ul>
57        <li><a href="/e/action/ListInfo/?classid=46">艺术设计系</a></li>
58        <li><a href="/e/action/ListInfo/?classid=47">机械电子系</a></li>
59        <li><a href="/e/action/ListInfo/?classid=49">财会信息系</a></li>
60        <li><a href="/e/action/ListInfo/?classid=50">商贸旅游系</a></li>
61        <li><a href="/e/action/ListInfo/?classid=51">建筑工程系</a></li>
62        <li><a href="/e/action/ListInfo/?classid=52">新疆学部</a></li>
63      </ul>
64    </li>
65    <li><a href="javascript:void(0)"><span style="color: #ff0000">招生宣
   ↪    传</span></a>
66      <ul>
67        <li><a href="/html/recruit/zsbm.html">招生报名</a></li>
68        <li><a href="/html/recruit/plan.html">招生简章</a></li>
69        <li><a href="/html/recruit/pro.html">专业介绍</a></li>
70        <li><a href="/html/recruit/faq.html">热点问答</a></li>
71        <li><a href="/html/recruit/xysh.html">校园生活</a></li>
```

```
72      </ul>
73    </li>
74    <li><a href="javascript:void(0)"><span>校务公开</span></a>
75      <ul>
76        <li><a href="/html/public/org.html">组织架构</a></li>
77        <li><a href="/e/action/ListInfo/?classid=19">办学规划</a></li>
78        <li><a href="/e/action/ListInfo/?classid=20">管理制度</a></li>
79        <li><a href="/e/action/ListInfo/?classid=21">阳光收费</a></li>
80        <li><a href="/e/action/ListInfo/?classid=22">评职评优</a></li>
81        <li><a href="/e/action/ListInfo/?classid=23">招标公告</a></li>
82        <li><a href="/e/action/ListInfo/?classid=390">质量报告</a></li>
83      </ul>
84    </li>
85    <li><a href="javascript:void(0)"><span>专题网站</span></a>
86      <ul>
87        <li><a href="/html/exemplary/about.html"><span>示范校专题
        ↪   网</span></a></li>
88        <li><a href="/htmlhomepage/yiheliangyi/about.html">一核二翼专题
        ↪   网</a></li>
89        <li><a href="/e/action/ListInfo/?classid=364">群众路线活动</a></li>
90        <li><a href="/e/action/ListInfo/?classid=371">旅游职教集团</a></li>
91        <li><a href="/e/action/ListInfo/?classid=39">德 育 品 牌</a></li>
92        <li><a href="http://server2.sxszjzx.com/~jwc">精 品 课 程</a></li>
93        <li><a href="/e/action/ListInfo/?classid=379">信 息 中 心</a></li>
94      </ul>
95    </li>
96    <li><a target="_blank" href="/e/action/ListInfo/?classid=387"><span
    ↪   style="color: #ff0000">党建工作</span></a></li>
97
98    </ul>
99  </div>
```

### 4.1.2 content

```
1 (def selector-content [:div.page_1])
```

1. sidebar

   `http://www.sxti.zj.cn/html/school/about.html`

   ```
   1 (def selector-sidebar [:div.page_left :div.pleft_t3])
   ```

   ```html
   1  <div class="page_left">
   2    ..
   3
   4    <div class="pleft_t2">
   5      <ul class="pleft_t3">
   6
   7        <li><a href="/html/school/about.html">学校简介</a></li>
   8        <li><a href="/e/action/ListInfo/?classid=13">现任领导</a></li>
   9        <li><a href="/e/action/ListInfo/?classid=283">名优教师</a></li>
   10       <li><a href="/html/school/office.html">处室联系</a></li>
   11       <li><a href="/e/action/ListInfo/?classid=14">校园风景</a></li>
   12       <li><a href="/e/action/ListInfo/?classid=35">楼层分布</a></li>
   13       <li><a href="/html/school/map.html">交通地图</a></li>
   14       <li><a href="/e/action/ListInfo/?classid=15">学校荣誉</a></li>
   15       <li><a href="/e/action/ListInfo/?classid=377">媒体关注</a></li>
   16
   17     </ul>
   18   </div>
   19
   20 </div>
   ```

2. article

```
1 (def selector-article [:div.page_right :div.pright_t3])
```

```
1 <div class="page_right">
2   <!-- title -->
3   <div class="pright_t3">
4     <!-- article -->
5     <div class="pright_t4">
6
7     </div>
8   </div>
9 </div>
```

## 4.2   **DONE** website-new.clj HTML elements

**CLOSED:** *[2018-11-19 Mon 19:11]*

### 4.2.1   nav

http://www.sxszjzx.com/

```
1 (def selector-nav [:div.nav])
```

```
1 <div class="nav">
2   <div class="siteWidth">
3
4     <ul id="mainNav" class="mainNav">
```

```html
<li class="li1 first1 on1" id="li-home">
  <h3 class="h1">
    <a class="a1" href="/">网站首页</a>
  </h3>
</li>

<li class="li1 hasUl1" id="li-xygk">
  <h3 class="h1">
    <a class="a1" href="/xygk/xyjj">学院概况</a></h3>
  <ul class="ul1" style="display: none;">

    <li class="li2 first2">
      <h3 class="h2"><a class="a2" href="/xygk/xyjj">学院简介</a></h3>
    </li>
    <li class="li2">
      <h3 class="h2"><a class="a2" href="/xygk/xrld">现任领导</a></h3>
    </li>
    <li class="li2">
      <h3 class="h2"><a class="a2" href="/xygk/zzjg">组织架构</a></h3>
    </li>
    <li class="li2">
      <h3 class="h2"><a class="a2" href="/xygk/cslx">处室联系</a></h3>
    </li>
    <li class="li2">
      <h3 class="h2"><a class="a2" href="/xygk/xyfg">校园风光</a></h3>
    </li>
    <li class="li2">
      <h3 class="h2"><a class="a2" href="/xygk/xyry">学院荣誉</a></h3>
    </li>
    <li class="li2">
      <h3 class="h2"><a class="a2" href="/xygk/lsyg">历史沿革</a></h3>
```

```
36        </li>
37        <li class="li2">
38          <h3 class="h2"><a class="a2" href="/xygk/lcfb">楼层分布</a></h3>
39        </li>
40        <li class="li2 last2">
41          <h3 class="h2"><a class="a2" href="/xygk/jtdt">交通地图</a></h3>
42        </li>

43
44      </ul>
45    </li>
46    <li class="li1 hasUl1" id="li-xydt">
47      <h3 class="h1">
48        <a class="a1" href="/xydt">学院动态</a></h3>
49      <ul class="ul1" style="display: none;">

50
51        <li class="li2 first2">
52          <h3 class="h2"><a class="a2" href="/xydt/xyxw"
            ↪   target="_blank">学院新闻</a></h3>
53        </li>
54        <li class="li2">
55          <h3 class="h2"><a class="a2" href="/xydt/xbxw">系部新闻</a></h3>
56        </li>
57        <li class="li2">
58          <h3 class="h2"><a class="a2" href="/xydt/mtjj">媒体聚焦</a></h3>
59        </li>
60        <li class="li2">
61          <h3 class="h2"><a class="a2" href="/xydt/xyry1">学院荣誉</a></h3>
62        </li>
63        <li class="li2">
64          <h3 class="h2"><a class="a2" href="/xydt/jsry">教师荣誉</a></h3>
65        </li>
```

```
66    <li class="li2 last2">
67      <h3 class="h2"><a class="a2" href="/xydt/xsry">学生荣誉</a></h3>
68    </li>
69
70    </ul>
71  </li>
72  <li class="li1 hasUl1" id="li-xbjs">
73    <h3 class="h1">
74      <a class="a1" href="/xbjs">系部建设</a></h3>
75    <ul class="ul1" style="display: none;">
76
77      <li class="li2 first2 hasUl2">
78        <h3 class="h2"><a class="a2" href="/xbjs/yssjx">艺术设计
        ↪   系</a></h3>
79      </li>
80      <li class="li2 hasUl2">
81        <h3 class="h2"><a class="a2" href="/xbjs/jxdzx">机械电子
        ↪   系</a></h3>
82      </li>
83      <li class="li2 hasUl2">
84        <h3 class="h2"><a class="a2" href="/xbjs/chxxx">财会信息
        ↪   系</a></h3>
85      </li>
86      <li class="li2 hasUl2">
87        <h3 class="h2"><a class="a2" href="/xbjs/smlyx">商贸旅游
        ↪   系</a></h3>
88      </li>
89      <li class="li2 hasUl2">
90        <h3 class="h2"><a class="a2" href="/xbjs/jzgcx">建筑工程
        ↪   系</a></h3>
91      </li>
```

```
92        <li class="li2 last2 hasUl2">
93          <h3 class="h2"><a class="a2" href="/xbjs/xjxb">新疆学部</a></h3>
94        </li>
95
96      </ul>
97    </li>
98    <li class="li1 hasUl1" id="li-zsjy">
99      <h3 class="h1">
100        <a class="a1" href="/zsjy/zsbm">招生就业</a></h3>
101      <ul class="ul1" style="display: none;">
102
103        <li class="li2 first2">
104          <h3 class="h2"><a class="a2" href="/zsjy/zsbm">招生报名</a></h3>
105        </li>
106        <li class="li2">
107          <h3 class="h2"><a class="a2" href="/zsjy/zsjz">招生简章</a></h3>
108        </li>
109        <li class="li2">
110          <h3 class="h2"><a class="a2" href="/zsjy/zyjs6">专业介绍</a></h3>
111        </li>
112        <li class="li2">
113          <h3 class="h2"><a class="a2" href="/zsjy/rdwd">热点问答</a></h3>
114        </li>
115        <li class="li2">
116          <h3 class="h2"><a class="a2" href="/zsjy/xysh">校园生活</a></h3>
117        </li>
118        <li class="li2">
119          <h3 class="h2"><a class="a2" href="/zsjy/zxbm">在线报名</a></h3>
120        </li>
121        <li class="li2 last2">
122          <h3 class="h2"><a class="a2" href="/zsjy/jyxx">就业信息</a></h3>
```

```
123        </li>

124

125      </ul>

126    </li>

127    <li class="li1 hasUl1" id="li-ztlm">

128      <h3 class="h1">

129        <a class="a1" href="/ztlm">专题栏目</a></h3>

130      <ul class="ul1" style="display: none;">

131

132        <li class="li2 first2 hasUl2">

133          <h3 class="h2"><a class="a2" href="/ztlm/sfxjs">示范校建

                ↪    设</a></h3>

134        </li>

135        <li class="li2 hasUl2">

136          <h3 class="h2"><a class="a2" href="/ztlm/smgc">三名工程</a></h3>

137        </li>

138        <li class="li2 hasUl2">

139          <h3 class="h2"><a class="a2" href="/ztlm/yheyzt">一核二翼专

                ↪    题</a></h3>

140        </li>

141        <li class="li2 hasUl2">

142          <h3 class="h2"><a class="a2" href="/ztlm/qzlxhd">群众路线活

                ↪    动</a></h3>

143        </li>

144        <li class="li2 hasUl2">

145          <h3 class="h2"><a class="a2" href="/ztlm/lyzjjt">旅游职教集

                ↪    团</a></h3>

146        </li>

147        <li class="li2 hasUl2">

148          <h3 class="h2"><a class="a2" href="/ztlm/dypp">德育品牌</a></h3>

149        </li>
```

```
150    <li class="li2 hasUl2">
151      <h3 class="h2"><a class="a2" href="/ztlm/xysp">校园视频</a></h3>
152    </li>
153    <li class="li2 last2">
154      <h3 class="h2"><a class="a2"
       ↪   href="http://server2.sxszjzx.com/~jwc/" target="_blank">精品
       ↪   课程</a></h3>
155    </li>
156
157  </ul>
158  </li>
159  <li class="li1 hasUl1" id="li-xxgk">
160    <h3 class="h1">
161      <a class="a1" href="/xxgk">校务公开</a></h3>
162    <ul class="ul1" style="display: none;">
163
164      <li class="li2 first2">
165        <h3 class="h2"><a class="a2" href="/xxgk/bxgh">办学规划</a></h3>
166      </li>
167      <li class="li2">
168        <h3 class="h2"><a class="a2" href="/xxgk/glzd">公示公告</a></h3>
169      </li>
170      <li class="li2">
171        <h3 class="h2"><a class="a2" href="/xxgk/ygsf">阳光收费</a></h3>
172      </li>
173      <li class="li2">
174        <h3 class="h2"><a class="a2" href="/xxgk/pzpy">评职评优</a></h3>
175      </li>
176      <li class="li2">
177        <h3 class="h2"><a class="a2" href="/xxgk/zbgg">招标公告</a></h3>
178      </li>
```

```
179        <li class="li2">
180          <h3 class="h2"><a class="a2" href="/xxgk/zlbg">质量报告</a></h3>
181        </li>
182        <li class="li2 last2">
183          <h3 class="h2"><a class="a2" href="/xxgk/zyxz">资源下载</a></h3>
184        </li>
185
186      </ul>
187    </li>
188    <li class="li1 last1 hasUl1" id="li-djgz">
189      <h3 class="h1">
190        <a class="a1" href="/djgz">党建工作</a></h3>
191      <ul class="ul1" style="display: none;">
192
193        <li class="li2 first2">
194          <h3 class="h2"><a class="a2" href="/djgz/djdt">党建动态</a></h3>
195        </li>
196        <li class="li2">
197          <h3 class="h2"><a class="a2" href="/djgz/lzzl">廉政专栏</a></h3>
198        </li>
199        <li class="li2">
200          <h3 class="h2"><a class="a2" href="/djgz/lqhd">亮旗行动</a></h3>
201        </li>
202        <li class="li2 last2">
203          <h3 class="h2"><a class="a2" href="/djgz/xxzl">学习资料</a></h3>
204        </li>
205
206      </ul>
207    </li>
208  </ul>
209
```

```
210  <script type="text/javascript">
211    (function () {
212    var navST;
213    var navST1;
214    var name = 'mainNav';
215    var t = 200;
216    var type = 1;
217    var removeOn = 'False';
218    var effect = 'slideDown';
219    var appendItem = '#';
220    var li = "#" + name + " li";
221    var index = 0;
222
223    if (!$("#" + name + " .li1").hasClass("on1")) {
224    $("#" + name + " .li1").first().addClass("on1");
225    } //默认第一个加.on1 类
226    index = $("#" + name + " .li1").index($("#" + name + " .on1"));
227
228    //鼠标离开导航后，回复默认.on1 类位置
229    $("#" + name)
230    .hover(
231    function () {
232    if (navST1 != null) {
233    clearTimeout(navST1);
234    }
235    },
236    function () {
237    navST1 = setTimeout(function () {
238    $("#" + name + " .li1").removeClass("on1").eq(index).addClass("on1");
239    },
240    500);
```

16

```javascript
241        }
242        );
243
244        if (type == '1') {
245        li = "#" + name + " .li1";
246        }
247        if (appendItem != '#') { //插入内容
248        var appendHtml = $(appendItem).html();
249        $(li).first().append(appendHtml);
250        $(appendItem).remove();
251        }
252
253        if (type == '3') {
254        $("#" + name + " .on1").find("ul").first().show();
255        }
256
257        $(li)
258        .hover(function () {
259        var curItem = $(this);
260        var onNum = (curItem.attr("class").split(" "))[0].replace("li", "");
261        $(li).removeClass("on" + onNum);
262        curItem.addClass("on" + onNum);
263        navST = setTimeout(function () { //延时触发
264
265        if ($("ul:first", curItem).css("display") != "block") {
266        $(li + " .ul" + onNum).hide();
267        if (effect == 'fade') {
268        $("ul:first", curItem).fadeIn(t);
269        } else {
270        $("ul:first", curItem).slideDown(t);
271        }
```

```
272    };
273    navST = null;
274    },
275    t);
276    },
277    function () {
278    if (navST != null) {
279    clearTimeout(navST);
280    }
281    if (type == '1' || type == '2') {
282    if (effect == 'fade') {
283    $(this).find("ul").first().fadeOut(t);
284    } else {
285    $(this).find("ul").first().slideUp(t);
286    }
287    }
288    if (removeOn == 'True') {
289    $(this).removeClass("on1");
290    }
291    },
292    t); //end hover
293    })()
294  </script>
295
296 </div>
297 </div>
```

### 4.2.2 content

```
1 (def selector-content [:div#content])
```

```
1 <div id="content">
2   ....
3 </div>
```

1. sidebar

   `http://www.sxszjzx.com/xygk/xyjj`

```clojure
1 (def selector-sidebar [:aside.side])
```

```
1 <aside class="side">
2
3   <div id="sideMenu">
4     <div class="hd">
5       <h3>学院概况</h3>
6     </div>
7     <div class="bd">
8       <ul class="menuList">
9         <li class="on"><a href="/xygk/xyjj">学院简介</a></li>
10        <li><a href="/xygk/xrld">现任领导</a></li>
11        <li><a href="/xygk/zzjg">组织架构</a></li>
12        <li><a href="/xygk/cslx">处室联系</a></li>
13        <li><a href="/xygk/xyfg">校园风光</a></li>
14        <li><a href="/xygk/xyry">学院荣誉</a></li>
15        <li><a href="/xygk/lsyg">历史沿革</a></li>
16        <li><a href="/xygk/lcfb">楼层分布</a></li>
17        <li><a href="/xygk/jtdt">交通地图</a></li>
18      </ul>
19    </div>
20  </div>
```

```
21
22    <div id="sideRmph" class="sideBox">
23      <div class="hd">
24        <h3>热点资讯</h3>
25      </div>
26      <div class="bd">
27        <ul class="infoListB">
28
29          <li class="noData">暂无资料</li>
30        </ul>
31      </div>
32    </div>
33
34  </aside>
```

2. article

`http://www.sxszjzx.com/xygk/xyjj`

```
1  (def selector-article [:div.mainContent])
```

```
1  <div class="mainContent">
2
3    <div class="mHd">
4      <div class="path">
5
6        <em>您的位置：</em><a href="/">首页</a>
7        &gt;<a href="/xygk/xyjj">学院概况</a>&gt;<a href="/xygk/xyjj">学
   ↪    院简介</a></div>
8      <h3>学院简介</h3>
```

```html
</div>
<div class="mBd">
  <!-- 正文内容 S -->
  <div class="articleCon">
    <div class="printArea" data-power-area="content">
      <!-- 标题 -->
      <h3 class="title">学院简介</h3>
      <div class="property">
        <span>【字体: <a href="javascript:;"
          data-power-command="reducefont">小</a> <a
          href="javascript:;"
          data-power-command="enlargefont">大</a>】</span>
      </div>
      <!-- 正文 -->
      <div class="conTxt" data-power-defaultfontsize="16"
        data-power-defaultlineheight="2"
        data-power-imgmaxwidth="800">
```

21

```
21    <div><strong style="padding: 0px; margin: 0px; outline:
↪     none; color: rgb(51, 51, 51); font-family:
↪     &quot;Microsoft Yahei&quot;; font-size: 14px;
↪     white-space: normal; background-color: rgb(255, 255,
↪     255)">绍兴技师学院（筹）绍兴市职教中心</strong><span
↪     style="color: #333333; font-family: &quot;Microsoft
↪     Yahei&quot;; font-size: 14px; background-color:
↪     #FFFFFF">创办于</span><strong style="padding: 0px;
↪     margin: 0px; outline: none; color: rgb(51, 51, 51);
↪     font-family: &quot;Microsoft Yahei&quot;; font-size:
↪     14px; white-space: normal; background-color: rgb(255,
↪     255, 255)">1958</strong><span style="color: #333333;
↪     font-family: &quot;Microsoft Yahei&quot;; font-size:
↪     14px; background-color: #FFFFFF">年，是以培养现代服务业和
↪     先进制造业技能人才为主，集学历教育、职业培训、技能鉴定为一
↪     体的综合性职业学校。</span><strong style="padding: 0px;
↪     margin: 0px; outline: none; color: rgb(51, 51, 51);
↪     font-family: &quot;Microsoft Yahei&quot;; font-size:
↪     14px; white-space: normal; background-color: rgb(255,
↪     255, 255)">1996</strong><span style="color: #333333;
↪     font-family: &quot;Microsoft Yahei&quot;; font-size:
↪     14px; background-color: #FFFFFF">年被评为首批国家级重点职
↪     业学校，</span><strong style="padding: 0px; margin: 0px;
↪     outline: none; color: rgb(51, 51, 51); font-family:
↪     &quot;Microsoft Yahei&quot;; font-size: 14px;
↪     white-space: normal; background-color: rgb(255, 255,
↪     255)">2013</strong><span style="color: #333333;
↪     font-family: &quot;Microsoft Yahei&quot;; font-size:
↪     14px; background-color: #FFFFFF">年被确定为首批国家中等职
↪     业教育改革发展示范学校。</span><br style="padding: 0px;
↪     margin: 0px; outline: none; color: rgb(51, 51, 51);
↪     font-family: &quot;Microsoft Yahei&quot;; font-size:
↪     14px; white-space: normal; background-color: rgb(255,
↪     255, 255)"><span style="color: #333333; font-family:
↪     &quot;Microsoft Yahei&quot;; font-size: 14px;
↪     background-color: #FFFFFF">     学校实行
↪     "</span><strong style="padding: 0px; margin: 0px;
↪     outline: none; color: rgb(51, 51, 51); font-family:
↪     &quot;Microsoft Yahei&quot;; font-size: 14px;
↪     white-space: normal; background-color: rgb(255, 255,
```

```html
22          </div>
23        </div>
24      <div class="userControl">
25
26
27        <div class="bdsharebuttonbox"><a href="#" class="bds_more"
   ↪   data-cmd="more"></a><a href="#" class="bds_qzone"
   ↪   data-cmd="qzone" title=" 分享到 QQ 空间"></a><a href="#"
   ↪   class="bds_tsina" data-cmd="tsina" title=" 分享到新浪微
   ↪   博"></a><a href="#" class="bds_tqq" data-cmd="tqq" title="
   ↪   分享到腾讯微博"></a><a href="#" class="bds_renren"
   ↪   data-cmd="renren" title=" 分享到人人网"></a><a href="#"
   ↪   class="bds_weixin" data-cmd="weixin" title=" 分享到微
   ↪   信"></a></div>
28      <script>
29        window._bd_share_config = {
30        "common": {
31        "bdSnsKey": {},
32        "bdText": "",
33        "bdMini": "2",
34        "bdMiniList": false,
35        "bdPic": "",
36        "bdStyle": "1",
37        "bdSize": "24"
38        },
39        "share": {}
40        };
41        with (document) {
42        0[(getElementsByTagName('head')[0] || body)
43        .appendChild(createElement('script'))
```

```
44            .src =
      ↪    'http://bdimg.share.baidu.com/static/api/js/share.js?v=89860593.js?cdnve
      ↪    +
45         ~(-new Date() / 36e5)];
46       }
47     </script>
48   </div>
49 </div>
50 <!-- 正文内容 E -->
51 </div>
52
53 </div>
```

# 5   TODO crawler.clj 获取所有文章的方法 [15/17]

## 5.1   DONE benchmarking the HTTP URL request speed [2/2]

**CLOSED:** *[2018-11-19 Mon 18:04]*

☒ record to Org Enlive.org

☒ optimize code in here

### 5.1.1   clj-http + enlive/html-snippet

```
1 (require '[clj-http.client :as http])
2 (require '[net.cgrand.enlive-html :as html])
3 (use 'criterium.core)
4
5 (pr (quick-bench
6     (-> (http/get "https://www.baidu.com")
```

```
7            :body
8            html/html-snippet)))
```

```
Evaluation count : 6 in 6 samples of 1 calls.
             Execution time mean : 260.639936 ms
    Execution time std-deviation : 74.939971 ms
   Execution time lower quantile : 161.264308 ms ( 2.5%)
   Execution time upper quantile : 337.518891 ms (97.5%)
                   Overhead used : 36.702468 ns
nil
```

### 5.1.2   enlive/html-resource + URL

```
1 (require '[net.cgrand.enlive-html :as html])
2 (import 'java.net.URL)
3 (use 'criterium.core)
4
5 (pr (quick-bench
6       (html/html-resource (URL. "https://www.baidu.com"))))
```

```
Evaluation count : 18 in 6 samples of 3 calls.
             Execution time mean : 64.137600 ms
    Execution time std-deviation : 13.476096 ms
   Execution time lower quantile : 49.122882 ms ( 2.5%)
   Execution time upper quantile : 79.881576 ms (97.5%)
                   Overhead used : 36.702468 ns
nil
```

## 5.2   DONE 按照网站的导航栏目录入口获取相应类目下的所有文章 [9/9]

**CLOSED:** *[2018-11-20 Tue 13:40]*

### 5.2.1 **DONE** old

**CLOSED:** *[2018-11-20 Tue 20:13]*

1. **DONE** get navigator bar links **CLOSED:** *[2018-11-20 Tue 10:50]*
   http://www.sxszjzx.com/html/school/about.html

```
1  <body>
2    <div class="page_all">
3      <div class="head_1">
4        <div class="head_2">
5          <div class="head_4">
6            <div class="page_1">
7              <div class="page_left">
8                <div class="page_right">
9                  <div class="pright_t3">
10                   <div class="pright_t4">
```

```
1  (def website-old-url "http://www.sxti.zj.cn")
2  (def website-old-html (get-html "http://www.sxti.zj.cn/"))
3
4  (defn get-html
5    "Get HTML string as result."
6    [url]
7    (-> (http/get url {:as "GB2312"})
8        :body
9        html/html-snippet))
10
11 (defonce nav-bar
12   (html/select
```

```
13      (drop 1
14           (first (map #(html/select % [:li])
15                       ;; nav bar
16                       (html/select
17                        website-old-html
18                        [:html :body :div.page_all :div.head_2
                          ↪  :div.head_3 :ul#head_nav]))))
19      [:a]))
20

21 (defonce nav-bar-links-map
22   (map #(let [link  (str website-old-url
23                         ;; :attrs nil (:href does not exist)
24                         (if (nil? (:attrs %))
25                           nil
26                           ;; :href "javascript:void(0)"
27                           (if (= (first (html/attr-values % :href))
                             ↪  "javascript:void(0)")
28                             nil
29                             ;; :href "/..."
30                             (first (html/attr-values % :href)))))
31                title (html/text %)]
32           {title link})
33         nav-bar))
34

35 (pprint nav-bar-links-map)
```

---

```
1 ({" 学校概况" "http://www.sxti.zj.cn"}
2  {" 学校简介" "http://www.sxti.zj.cn/html/school/about.html"}
3  {" 现任领导" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=13"}
4  {" 名优教师" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=283"}
```

```
5   {" 处室联系" "http://www.sxti.zj.cn/html/school/office.html"}
6   {" 校园风景" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=14"}
7   {" 楼层分布" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=35"}
8   {" 交通地图" "http://www.sxti.zj.cn/html/school/map.html"}
9   {" 学校荣誉" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=15"}
10  {" 媒体关注" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=377"}
11  {" 学校简介" "http://www.sxti.zj.cn/html/school/about.html"}
12  {" 现任领导" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=13"}
13  {" 名优教师" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=283"}
14  {" 处室联系" "http://www.sxti.zj.cn/html/school/office.html"}
15  {" 校园风景" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=14"}
16  {" 楼层分布" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=35"}
17  {" 交通地图" "http://www.sxti.zj.cn/html/school/map.html"}
18  {" 学校荣誉" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=15"}
19  {" 媒体关注" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=377"}
20  {" 新闻频道" "http://www.sxti.zj.cn"}
21  {" 校园新闻" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=33"}
22  {" 系部新闻" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=34"}
23  {" 对外交流" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=7"}
24  {" 学校荣誉" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=8"}
25  {" 师生荣誉" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=9"}
26  {" 校园视频" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=10"}
27  {" 二周安排" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=11"}
28  {" 每周寄语" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=12"}
29  {" 校园新闻" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=33"}
30  {" 系部新闻" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=34"}
31  {" 对外交流" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=7"}
32  {" 学校荣誉" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=8"}
33  {" 师生荣誉" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=9"}
34  {" 校园视频" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=10"}
35  {" 二周安排" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=11"}
```

```
36  {" 每周寄语" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=12"}
37  {" 下载频道" "http://www.sxti.zj.cn"}
38  {" 党政办" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=16"}
39  {" 教务处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=17"}
40  {" 德育团委" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=18"}
41  {" 总务处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=53"}
42  {" 实习处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=54"}
43  {" 安保处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=56"}
44  {" 教科室" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=175"}
45  {" 教学软件" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=57"}
46  {" 其它" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=58"}
47  {" 党政办" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=16"}
48  {" 教务处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=17"}
49  {" 德育团委" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=18"}
50  {" 总务处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=53"}
51  {" 实习处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=54"}
52  {" 安保处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=56"}
53  {" 教科室" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=175"}
54  {" 教学软件" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=57"}
55  {" 其它" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=58"}
56  {" 处室网站" "http://www.sxti.zj.cn"}
57  {" 党政办" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=37"}
58  {" 教务处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=38"}
59  {" 德育团委" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=39"}
60  {" 总务处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=40"}
61  {" 实习处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=41"}
62  {" 安保处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=43"}
63  {" 教科室" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=104"}
64  {" 服务公司" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=44"}
65  {" 党政办" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=37"}
66  {" 教务处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=38"}
```

```
67    {" 德育团委" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=39"}
68    {" 总务处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=40"}
69    {" 实习处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=41"}
70    {" 安保处" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=43"}
71    {" 教科室" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=104"}
72    {" 服务公司" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=44"}
73    {" 系部网站" "http://www.sxti.zj.cn"}
74    {" 艺术设计系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=46"}
75    {" 机械电子系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=47"}
76    {" 财会信息系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=49"}
77    {" 商贸旅游系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=50"}
78    {" 建筑工程系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=51"}
79    {" 新疆学部" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=52"}
80    {" 艺术设计系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=46"}
81    {" 机械电子系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=47"}
82    {" 财会信息系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=49"}
83    {" 商贸旅游系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=50"}
84    {" 建筑工程系" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=51"}
85    {" 新疆学部" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=52"}
86    {" 招生宣传" "http://www.sxti.zj.cn"}
87    {" 招生报名" "http://www.sxti.zj.cn/html/recruit/zsbm.html"}
88    {" 招生简章" "http://www.sxti.zj.cn/html/recruit/plan.html"}
89    {" 专业介绍" "http://www.sxti.zj.cn/html/recruit/pro.html"}
90    {" 热点问答" "http://www.sxti.zj.cn/html/recruit/faq.html"}
91    {" 校园生活" "http://www.sxti.zj.cn/html/recruit/xysh.html"}
92    {" 招生报名" "http://www.sxti.zj.cn/html/recruit/zsbm.html"}
93    {" 招生简章" "http://www.sxti.zj.cn/html/recruit/plan.html"}
94    {" 专业介绍" "http://www.sxti.zj.cn/html/recruit/pro.html"}
95    {" 热点问答" "http://www.sxti.zj.cn/html/recruit/faq.html"}
96    {" 校园生活" "http://www.sxti.zj.cn/html/recruit/xysh.html"}
97    {" 校务公开" "http://www.sxti.zj.cn"}
```

```
98  {" 组织架构" "http://www.sxti.zj.cn/html/public/org.html"}
99  {" 办学规划" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=19"}
100  {" 管理制度" "http://www.sxti.zj.cn/e/action/ListInfo/?classid=20"}
101  ...)
```

2. **DONE** parse sub-nav page to get all listed result articles **CLOSED:** *[2018-11-20 Tue 19:56]*

```
1  <a
   ↪  href="/e/action/ListInfo/index.php?page=109&amp;classid=33&amp;totalnum=1648">尾
   ↪  页</a>
```

http://www.sxti.zj.cn/e/action/ListInfo/index.php?page=15&classid=33&totalnum=1648

```
1  (defn- get-mainContent-html
2    "Get the <div class=\"page_right\" element."
3    [html]
4    (html/select html [:html :body :div.page_all :div.page_1
     ↪   :div.page_right]))
5
6  (map #(let [link  (str website-old-url (first (html/attr-values %
   ↪  :href)))
7              title (html/text %)]
8          {title link})
9      (html/select
10       (get-mainContent-html
11        (get-html "http://www.sxti.zj.cn/e/action/ListInfo/?classid=33"
          ↪   {:as "GB2312"}))
12       [:ul.newsList1 :li :a]))
```

```
13

14

15 (defn- get-mainContent-html
16   "Get the <div class=\"page_right\" element."
17   [html]
18   (html/select html [:html :body :div.page_all :div.page_1
     ↪  :div.page_right]))

19

20 (defn get-page-article-links
21   "Get articles list's every article link and title."
22   [nav-link]
23   (map #(let [link  (str website-old-url (first (html/attr-values %
     ↪  :href)))
24              title (html/text %)]
25         {title link})
26      (html/select
27       (get-mainContent-html
28        (get-html nav-link {:as "GB2312"}))
29        [:ul.newsList1 :li :a])))
```

3. **DONE** iterate all result pages **CLOSED:** *[2018-11-20 Tue 20:13]*

```
1 (defn get-total-result-pages
2   "How much result pages?"
3   [nav-link]
4   (Integer.
5    ((keyword (str "/e/action/ListInfo/index.php?" "page"))
6     (clojure.walk/keywordize-keys
7      (ring.util.codec/form-decode
8       (first
```

```
 9        (html/attr-values
10         (last
11          (html/select
12           (get-mainContent-html (get-html nav-link {:as "GB2312"}))
13           [:div.yema1 :a]))
14         :href)))))))

15

16 (comment
17   (get-total-result-pages
   ↪   "http://www.sxti.zj.cn/e/action/ListInfo/?classid=33"))

18

19 (defn get-all-article-links
20   "Get a nav's all articles link and title with map as return."
21   [nav-link]
22   (for [n (range 1 (inc (get-total-result-pages nav-link)))]
23     (let [url (str nav-link "&page=" n)]
24       (get-page-article-links url))))

25

26 ;;
   ↪   "http://www.sxti.zj.cn/e/action/ListInfo/index.php?classid=33&page=1&totalnum=1648

27

28 (get-all-article-links
   ↪   "http://www.sxti.zj.cn/e/action/ListInfo/?classid=33")
```

### 5.2.2  DONE new

**CLOSED:** *[2018-11-20 Tue 16:50]*

- http://zjzx.sxsedu.net/

1. **DONE** get navigator bar links **CLOSED:** *[2018-11-20 Tue 13:40]*

- http://zjzx.sxsedu.net/xygk/xyjj

---

```html
<body>
  <header id="header">
    <div class="nav">
      <div id="content">
        <aside class="side">
          <div class="mainContent">
            <div class="mHd">
              <h3>学院简介</h3>
              <div class="mBd">
                <div class="articleCon">
                  <div class="printArea">
                    <h3 class="title">学院简介</h3>
                    <div class="conTxt">
                      <div>
                        <footer id="footer">
```

---

```clojure
(require '[clj-http.client :as http]
         '[net.cgrand.enlive-html :as html])
(import 'java.net.URL)

(defn get-html
  "Get HTML string as result."
  [url]
  ;; faster, use tagsoup internal. But I don't know how to specify
  ↪  encoding.
  (html/html-resource (URL. url)))

(def website-new-url "http://www.sxszjzx.com")
```

34

```clojure
12 (def website-new-html (get-html website-new-url))

13

14 ;;; Nav sections
15 (def nav-bar
16   (drop 1 (html/select
17             website-new-html
18             [:html :body :div.wrap :div.nav
19              :div.siteWidth :ul#mainNav.mainNav
20              :li])))

21

22 (def nav-links
23   (map #(let [nav   (-> (html/select % [:a])
24                          first)
25               link  (str website-new-url (first (html/attr-values nav
                ↪  :href)))
26               title (html/text nav)]
27           {title link})
28        nav-bar))

29

30 (pprint nav-links)
```

```clojure
1 ({" 学院概况" "http://www.sxszjzx.com/xygk/xyjj"}
2  {" 学院简介" "http://www.sxszjzx.com/xygk/xyjj"}
3  {" 现任领导" "http://www.sxszjzx.com/xygk/xrld"}
4  {" 组织架构" "http://www.sxszjzx.com/xygk/zzjg"}
5  {" 处室联系" "http://www.sxszjzx.com/xygk/cslx"}
6  {" 校园风光" "http://www.sxszjzx.com/xygk/xyfg"}
7  {" 学院荣誉" "http://www.sxszjzx.com/xygk/xyry"}
8  {" 历史沿革" "http://www.sxszjzx.com/xygk/lsyg"}
9  {" 楼层分布" "http://www.sxszjzx.com/xygk/lcfb"}
```

```
10    {" 交通地图" "http://www.sxszjzx.com/xygk/jtdt"}
11    {" 学院动态" "http://www.sxszjzx.com/xydt"}
12    {" 学院新闻" "http://www.sxszjzx.com/xydt/xyxw"}
13    {" 系部新闻" "http://www.sxszjzx.com/xydt/xbxw"}
14    {" 媒体聚焦" "http://www.sxszjzx.com/xydt/mtjj"}
15    {" 学院荣誉" "http://www.sxszjzx.com/xydt/xyry1"}
16    {" 教师荣誉" "http://www.sxszjzx.com/xydt/jsry"}
17    {" 学生荣誉" "http://www.sxszjzx.com/xydt/xsry"}
18    {" 系部建设" "http://www.sxszjzx.com/xbjs"}
19    {" 艺术设计系" "http://www.sxszjzx.com/xbjs/yssjx"}
20    {" 机械电子系" "http://www.sxszjzx.com/xbjs/jxdzx"}
21    {" 财会信息系" "http://www.sxszjzx.com/xbjs/chxxx"}
22    {" 商贸旅游系" "http://www.sxszjzx.com/xbjs/smlyx"}
23    {" 建筑工程系" "http://www.sxszjzx.com/xbjs/jzgcx"}
24    {" 新疆学部" "http://www.sxszjzx.com/xbjs/xjxb"}
25    {" 招生就业" "http://www.sxszjzx.com/zsjy/zsbm"}
26    {" 招生报名" "http://www.sxszjzx.com/zsjy/zsbm"}
27    {" 招生简章" "http://www.sxszjzx.com/zsjy/zsjz"}
28    {" 专业介绍" "http://www.sxszjzx.com/zsjy/zyjs6"}
29    {" 热点问答" "http://www.sxszjzx.com/zsjy/rdwd"}
30    {" 校园生活" "http://www.sxszjzx.com/zsjy/xysh"}
31    {" 在线报名" "http://www.sxszjzx.com/zsjy/zxbm"}
32    {" 就业信息" "http://www.sxszjzx.com/zsjy/jyxx"}
33    {" 专题栏目" "http://www.sxszjzx.com/ztlm"}
34    {" 示范校建设" "http://www.sxszjzx.com/ztlm/sfxjs"}
35    {" 三名工程" "http://www.sxszjzx.com/ztlm/smgc"}
36    {" 一核二翼专题" "http://www.sxszjzx.com/ztlm/yheyzt"}
37    {" 群众路线活动" "http://www.sxszjzx.com/ztlm/qzlxhd"}
38    {" 旅游职教集团" "http://www.sxszjzx.com/ztlm/lyzjjt"}
39    {" 德育品牌" "http://www.sxszjzx.com/ztlm/dypp"}
40    {" 校园视频" "http://www.sxszjzx.com/ztlm/xysp"}
```

```
41  {" 精品课程"
    ↪  "http://www.sxszjzx.comhttp://server2.sxszjzx.com/~jwc/"}
42  {" 校务公开" "http://www.sxszjzx.com/xxgk"}
43  {" 办学规划" "http://www.sxszjzx.com/xxgk/bxgh"}
44  {" 公示公告" "http://www.sxszjzx.com/xxgk/glzd"}
45  {" 阳光收费" "http://www.sxszjzx.com/xxgk/ygsf"}
46  {" 评职评优" "http://www.sxszjzx.com/xxgk/pzpy"}
47  {" 招标公告" "http://www.sxszjzx.com/xxgk/zbgg"}
48  {" 质量报告" "http://www.sxszjzx.com/xxgk/zlbg"}
49  {" 资源下载" "http://www.sxszjzx.com/xxgk/zyxz"}
50  {" 党建工作" "http://www.sxszjzx.com/djgz"}
51  {" 党建动态" "http://www.sxszjzx.com/djgz/djdt"}
52  {" 廉政专栏" "http://www.sxszjzx.com/djgz/lzzl"}
53  {" 亮旗行动" "http://www.sxszjzx.com/djgz/lqhd"}
54  {" 学习资料" "http://www.sxszjzx.com/djgz/xxzl"})
```

2. **DONE** parse sub-nav page to get all listed result articles **CLOSED:** *[2018-11-20 Tue 15:44]*

```
1   <div class="mainContent">

2

3     <!-- nav name -->
4     <div class="mHd">
5       <div class="path">

6

7         <em>您的位置: </em><a href="/">首页</a>
8         &gt;<a href="/xydt">学院动态</a>&gt;<a href="/xydt/xyxw"
          ↪  target="_blank">学院新闻</a></div>
9       <h3>学院新闻</h3>
10    </div>
```

```
11
12  <div class="mBd">
13    <!-- 正文内容 S -->
14    <ul class="pageTPList">
15
16      <!-- article -->
17      <li class="first">
18
19        <div class="title">
20          <a target="_blank" class="tit"
          ↪   href="/xydt/xyxw/content_39935" title=" 六十载同心同德建
          ↪   名校  一甲子匠智匠力创品牌">六十载同心同德建名校  一甲子匠
          ↪   智匠力创品牌</a>
21        </div>
22
23        <div class="pic">
24          <a target="_blank" href="/xydt/xyxw/content_39935">
25            <img alt=" 六十载同心同德建名校  一甲子匠智匠力创品牌"
            ↪   src="/upload/sxszjzx/contentmanage/article/image/2018/11/09/20a47dec28
26          </a>
27        </div>
28
29        <div class="con">
30
31          <div class="intro">
32             薪火相承，一甲子春华秋实。2018 年，绍兴技师学院（筹）、
            ↪   绍兴市职业教育中心迎来了建校六十周年华诞。学校始建于
            ↪   1958 年，前身为鉴湖公社初级中学，1985 年更名为绍兴市第一
            ↪   职业技术中学，1995 年市一职中与市中兴职中和市二职中合并
            ↪   为绍兴市职教中心，1997 年市树人...
33          </div>
```

```html
34
        <div class="others">
35
          <span class="date">2018-11-09</span>
36
        </div>
37

38
      </div>
39

40
    </li>
41

42

43

44
    <li>
45

46
      <div class="title">
47
        <a target="_blank" class="tit"
48
        ↪  href="/xydt/xyxw/content_39927" title=" 你的样子，我刚好
        ↪  喜欢  ——我校开展"仪容仪表示范班"评比活动">你的样子，我
        ↪  刚好喜欢  ——我校开展"仪容仪表示范班"评比活动</a>
      </div>
49

50

51

52
      <div class="con">
53

54
        <div class="intro">
55
            为进一步加强学生的文明礼仪教育，强化学校的常规管
          ↪  理，使学生养成良好的行为习惯，促进优良的校风、班风的形成，
          ↪  近期德育处开展了"仪容仪表示范班"的评选活动。  根
          ↪  据学校的实际情况，学校制定了仪容仪表示范班的评选条件。要
          ↪  求男生不烫发染发，前面头发不盖...
        </div>
56

57
```

```html
58        <div class="others">
59          <span class="date">2018-11-08</span>
60        </div>

62      </div>

64    </li>



68    <li>

70      <div class="title">
71        <a target="_blank" class="tit"
     ↪    href="/xydt/xyxw/content_39923" title=" 当快闪遇上诗歌
     ↪    ——我们一起告白我的国">当快闪遇上诗歌——我们一起告白我的
     ↪    国</a>
72      </div>

74      <div class="pic">
75        <a target="_blank" href="/xydt/xyxw/content_39923">
76          <img alt=" 当快闪遇上诗歌——我们一起告白我的国"
     ↪      src="/upload/sxszjzx/contentmanage/article/image/2018/11/08/399e2baba1
77        </a>
78      </div>

80      <div class="con">
81
82        <div class="intro">
```

```
83                  近日，德育处、团委举办了一场"颂中华诗词，
   ↪      寻文化基因"的诗歌快闪活动。新颖的活动形式吸引了不少师生
   ↪      的关注，得到了许多同学的喝彩。"从浩瀚的地球仪上，我认识了
   ↪      我的祖国……"突然，熙熙攘攘的人群中传来了诗歌朗诵声，这
   ↪      边朗诵声初歇，那边人群中又突...
84          </div>
85
86          <div class="others">
87            <span class="date">2018-11-08</span>
88          </div>
89
90        </div>
91
92      </li>
93
94
95
96      <li>
97
98        <div class="title">
99          <a target="_blank" class="tit"
   ↪      href="/xydt/xyxw/content_39893" title=""红色匠心，青春向
   ↪      党"—我校十月"祖国颂"诗歌朗诵比赛圆满结束"> "红色匠心，
   ↪      青春向党"—我校十月"祖国颂"诗歌朗诵比赛圆满结束</a>
100       </div>
101
102       <div class="pic">
103         <a target="_blank" href="/xydt/xyxw/content_39893">
104           <img alt="  "红
   ↪      色匠心，青春向党"—我校十月"祖国颂"诗歌朗诵比赛圆满结束"
   ↪      src="/upload/sxszjzx/contentmanage/article/image/2018/11/07/1b861b6320
```

```
105            </a>
106        </div>
107
108        <div class="con">
109
110          <div class="intro">
111                为庆祝祖国 69 岁华诞，培养学生爱党爱国爱校情
              ↪  怀，配合省全民终身学习宣传周活动，提升学生诗歌朗诵水平及
              ↪  语文素养，营造朝气蓬勃、积极向上的校园文化氛围，11 月 1
              ↪  日下午，绍兴技师学院（筹）、绍兴市职教中心在报告厅隆重举行
              ↪  十月"祖国颂"诗歌朗诵比赛决赛。...
112          </div>
113
114          <div class="others">
115            <span class="date">2018-11-07</span>
116          </div>
117
118        </div>
119
120      </li>
121
122
123
124      <li>
125
126        <div class="title">
127          <a target="_blank" class="tit"
              ↪  href="/xydt/xyxw/content_39799" title=" 第三届校教职工气
              ↪  排球圆满结束 ">第三届校教职工气排球圆满结束 </a>
128        </div>
129
```

```html
130            <div class="pic">
131              <a target="_blank" href="/xydt/xyxw/content_39799">
132                <img alt=" 第三届校教职工气排球圆满结束 "
     ↪           src="/upload/sxszjzx/contentmanage/article/image/2018/11/02/bf87db2b58
133              </a>
134            </div>
135
136            <div class="con">
137
138              <div class="intro">
139                    第三届校教职工气排球比赛，经过二轮 26 场激烈比
     ↪           拚于周二中午在校体育馆落下帷幕。最终比赛成绩如下：第一名：
     ↪           新疆学部，第二名：办公行政教科组，第三名：德育实习安监组，
     ↪           第四名：艺术设计组。本届气排球比赛，由校教职工 10 个工会
     ↪           小组分别组队参加，比赛共分...
140              </div>
141
142              <div class="others">
143                <span class="date">2018-11-02</span>
144              </div>
145
146            </div>
147
148          </li>
149
150
151
152          <li>
153
154            <div class="title">
```

43

```html
155          <a target="_blank" class="tit"
    ↪    href="/xydt/xyxw/content_39808" title=" 关注课堂，携手育
    ↪    人——我校开展家长开放日观摩课活动报道">关注课堂，携手育人
    ↪    ——我校开展家长开放日观摩课活动报道</a>
156      </div>
157
158      <div class="pic">
159        <a target="_blank" href="/xydt/xyxw/content_39808">
160          <img alt=" 关注课堂，携手育人——我校开展家长开放日观摩课活
    ↪    动报道"
    ↪    src="/upload/sxszjzx/contentmanage/article/image/2018/11/04/6a62a56c68
161        </a>
162      </div>
163
164      <div class="con">
165
166        <div class="intro">
167          为更好的架设学校、家庭、社会沟通的桥梁，让家长们走进课堂，走
    ↪    近孩子，走进学校，10 月 24 日上午我校举行了"家长开放日"
    ↪    活动。"请家长听一堂课"作为这次活动的重头戏，学校和各系部
    ↪    都高度重视。除了语文、数学等文化课以外，各系部还展示了包
    ↪    括声乐、建筑设计、...
168        </div>
169
170        <div class="others">
171          <span class="date">2018-10-30</span>
172        </div>
173
174      </div>
175
176    </li>
```

```
177
178
179
180        <li>
181
182            <div class="title">
183              <a target="_blank" class="tit"
              ↪   href="/xydt/xyxw/content_39589" title=" "更高、更快、更
              ↪   强" ——绍兴技师学院（筹）绍兴市职教中心成功举办第七十一届
              ↪   田径运动会"> "更高、更快、更强" ——绍兴技师学院（筹）绍兴
              ↪   市职教中心成功举办第七十一届田径运动会</a>
184            </div>
185
186            <div class="pic">
187              <a target="_blank" href="/xydt/xyxw/content_39589">
188                <img alt=" "更高、更快、更强" ——绍兴技师学院（筹）绍兴市职
              ↪   教中心成功举办第七十一届田径运动会"
              ↪   src="/upload/sxszjzx/contentmanage/article/image/2018/10/28/4ea0ecdfcc
189              </a>
190            </div>
191
192            <div class="con">
193
194              <div class="intro">
195                 10 月 17 日至 19 日，我校隆重举办了第七十一届田径运动
              ↪   会。本届运动会共六个组别、十二个赛项，有 89 个班级、共
              ↪   1350 名运动员参赛，参赛班级、参赛人数均创历史新高。17 日
              ↪   下午，学校体育场上彩旗飘扬、音乐嘹亮，全校师生云集喜迎第
              ↪   七十一届田径运动会开幕。学校纪委书...
196              </div>
197
```

```
198            <div class="others">
199              <span class="date">2018-10-28</span>
200            </div>
201
202          </div>
203
204      </li>
205
206
207
208      <li>
209
210          <div class="title">
211            <a target="_blank" class="tit"
              ↪  href="/xydt/xyxw/content_39560" title=" 弘扬垦荒精神 铸
              ↪  牢党性之魂——我校党员教师赴大陈岛开展专题党日活动">弘扬垦
              ↪  荒精神 铸牢党性之魂——我校党员教师赴大陈岛开展专题党日活
              ↪  动</a>
212          </div>
213
214          <div class="pic">
215            <a target="_blank" href="/xydt/xyxw/content_39560">
216              <img alt=" 弘扬垦荒精神 铸牢党性之魂——我校党员教师赴大陈岛
                ↪  开展专题党日活动"
                ↪  src="/upload/sxszjzx/contentmanage/article/image/2018/10/26/f65137cdca
217            </a>
218          </div>
219
220          <div class="con">
221
222            <div class="intro">
```

```
223              为纪念中华人民共和国成立 69 周年，全面推进党的组织建设，激励
         ↪    全体党员教师继承和发扬党的优良传统和作风，绍兴市职教中心
         ↪    党委组织党员教师沿着习近平总书记的足迹奔赴浙江省直机关党
         ↪    员干部教育基地台州大陈岛开展现场学习教育，重温入党誓词，
         ↪    追溯红色记忆，学习 "...
224          </div>
225
226          <div class="others">
227            <span class="date">2018-10-26</span>
228          </div>
229
230        </div>
231
232      </li>
233
234
235
236      <li>
237
238        <div class="title">
239          <a target="_blank" class="tit"
         ↪    href="/xydt/xyxw/content_39499" title=" 加强师资建设，争
         ↪    创名师团队——我校召开《人性的追问与教师的职业成长》专题讲
         ↪    座">加强师资建设，争创名师团队——我校召开《人性的追问与教
         ↪    师的职业成长》专题讲座</a>
240        </div>
241
242        <div class="pic">
243          <a target="_blank" href="/xydt/xyxw/content_39499">
```

```html
244        <img alt=" 加强师资建设，争创名师团队——我校召开《人性的追
        ↪    问与教师的职业成长》专题讲座"
        ↪    src="/upload/sxszjzx/contentmanage/article/image/2018/10/24/cfc07cf774
245      </a>
246    </div>

247
248    <div class="con">

249
250      <div class="intro">
251        加强师资建设，争创名师团队——我校召开《人性的追问与教师的职
        ↪    业成长》专题讲座 10 月 12 日下午，学校邀请平湖职业中专校
        ↪    长贺陆军到我校报告厅召开《人性的追问与教师的职业成长》专
        ↪    题讲座。讲座由校长钱金星主持。全体教工认真聆听了讲座。贺
        ↪    校长从人的物质性与文化...
252      </div>

253
254      <div class="others">
255        <span class="date">2018-10-24</span>
256      </div>

257
258    </div>

259
260    </li>

261

262

263
264    <li class="last">

265
266      <div class="title">
```

```html
267        <a target="_blank" class="tit"
    ↪   href="/xydt/xyxw/content_39157" title=" 五星三名·用行动
    ↪   践行先锋力量——我校开展系列 "党员育人" 活动">五星三名·用
    ↪   行动践行先锋力量——我校开展系列 "党员育人" 活动</a>
268      </div>

269

270      <div class="pic">
271        <a target="_blank" href="/xydt/xyxw/content_39157">
272          <img alt=" 五
    ↪   星三名·用行动践行先锋力量——我校开展系列"党员育人"活动"
    ↪   src="/upload/sxszjzx/contentmanage/article/image/2018/10/10/0e90c18e22
273        </a>
274      </div>

275

276      <div class="con">

277

278        <div class="intro">
279          为充分发挥党员教师的先锋模范作用，进一步拓宽党员育人的广度与
    ↪   深度，传播先锋力量，本学期伊始，我校党委组织六大支部开展
    ↪   了 "秉烛怀志，躬身明责" 的系列党员育人岗活动—— "党员育
    ↪   人示范岗"，用红正的党徽为学生的放学之路保驾护航；"党建带
    ↪   团建"，用党员力...
280        </div>

281

282        <div class="others">
283          <span class="date">2018-09-29</span>
284        </div>

285

286      </div>

287

288    </li>
```

49

```
289
290    </ul>
291    <div class="page">
292

293

294      <a class="first disabled" href="javascript:void(0);">首页</a>
295      <a class="prev disabled" href="javascript:void(0);">上一页</a>
296      <a class="current">1</a> <a href="/xydt/xyxw_2">2</a> <a
    ↪   href="/xydt/xyxw_3">3</a> <a href="/xydt/xyxw_4">4</a> <a
    ↪   href="/xydt/xyxw_5">5</a> <a href="/xydt/xyxw_6">6</a> <a
    ↪   href="/xydt/xyxw_7">7</a> <a href="/xydt/xyxw_8">8</a> <a
    ↪   href="/xydt/xyxw_9">9</a>            <a class="next"
    ↪   href="/xydt/xyxw_2">下一页</a>
297      <a class="last" href="/xydt/xyxw_166">尾页</a>
298      <span class="total">共 1652 条信息/共 166 页</span>
299      <span class="select">转到第<input title=" 按回车键跳转到指定页"
    ↪   onkeypress="javascript:return quickJumpPage(event,
    ↪   'quickJumpButton')" type="text" style="width:24px"
    ↪   id="quickJumpInput" class="quickJumpInput "
    ↪   value="1"><script type="text/javascript">var pageNameUrl =
    ↪   '/xydt/xyxw_{pageid}';function quickJumpPage(event){if
    ↪   (event.keyCode == 13 && !(event.srcElement &&
    ↪   (event.srcElement.tagName.toLowerCase() == 'textarea')))
    ↪   {var number =
    ↪   document.getElementById('quickJumpInput').value;if
    ↪   (/\d+/i.test(number)){if (number > 166) number =
    ↪   166;if(number < 1) number = 1;window.location.href =
    ↪   pageNameUrl.replace("{pageid}", number);}else{alert(' 输入的
    ↪   页数有误！ ');}}}</script>页</span>
300    </div>
301    <!-- 正文内容 E -->
```

```
302    </div>
303
304  </div>
```

```clojure
1   '(defn get-mainContent-html
2      "Get nav link page's mainContent."
3      [nav-link]
4      (html/select nav-link [:html :body :div.wrap :div#content
       ↪   :div.mainContent]))
5
6   (defn get-articles-list
7      "Get page right articles list."
8      [nav-link]
9      (html/select (get-mainContent-html nav-link) [:div.mBd
       ↪   :ul.pageTPList
10                                                    :li :div.title
                                                      ↪   :a.tit]))
11
12  (pprint
13    (map #(let [link  (first (html/attr-values % :href))
14                title (html/text %)]
15            {title link})
16         (html/select
17          (html/select
18           (get-html "http://www.sxszjzx.com/xydt/xyxw")
19           [:html :body :div.wrap :div#content :div.mainContent])
20          [:div.mBd :ul.pageTPList :li :div.title :a.tit])))
```

```
1 ({" 六十载同心同德建名校  一甲子匠智匠力创品牌"
   ↪ "/xydt/xyxw/content_39935"}
2 {" 你的样子，我刚好喜欢  ——我校开展 "仪容仪表示范班" 评比活动"
   ↪ "/xydt/xyxw/content_39927"}
3 {" 当快闪遇上诗歌——我们一起告白我的国" "/xydt/xyxw/content_39923"}
4 {" "红色匠心，青春向党" —我校十月 "祖国颂" 诗歌朗诵比赛圆满结束"
   ↪ "/xydt/xyxw/content_39893"}
5 {" 第三届校教职工气排球圆满结束 " "/xydt/xyxw/content_39799"}
6 {" 关注课堂，携手育人——我校开展家长开放日观摩课活动报道"
   ↪ "/xydt/xyxw/content_39808"}
7 {" "更高、更快、更强" ——绍兴技师学院（筹）绍兴市职教中心成功举办第七十
   ↪ 一届田径运动会"
8 "/xydt/xyxw/content_39589"}
9 {" 弘扬垦荒精神 铸牢党性之魂——我校党员教师赴大陈岛开展专题党日活动"
   ↪ "/xydt/xyxw/content_39560"}
10 {" 加强师资建设，争创名师团队——我校召开《人性的追问与教师的职业成长》
   ↪ 专题讲座" "/xydt/xyxw/content_39499"}
11 {" 五星三名·用行动践行先锋力量——我校开展系列 "党员育人" 活动"
   ↪ "/xydt/xyxw/content_39157"})
```

3. **DONE** iterate all result pages **CLOSED:** *[2018-11-20 Tue 16:49]*

```
1 <span class="total">共 1652 条信息/共 166 页</span>
```

```
1 (defn get-total-result-pages
2   "How much result pages?"
3   [nav-link]
4   (Integer.
```

```clojure
5    (second
6     (re-find #"/共 (.*) 页"
7               (html/text
8                (first (html/select (get-mainContent-html (get-html
     ↪  nav-link))
9                                    [:div.mBd :div.page
                                     ↪  :span.total]))))))))
10
11 (pprint (get-total-result-pages "http://www.sxszjzx.com/xydt/xyxw"))
```

```
166
```

```
http://www.sxszjzx.com/xydt/xyxw
OR:
http://www.sxszjzx.com/xydt/xyxw_1

http://www.sxszjzx.com/xydt/xyxw_2
....
http://www.sxszjzx.com/xydt/xyxw_166
```

```clojure
1  (defn get-page-article-links
2    "Get articles list's every article link and title."
3    [nav-link]
4    (map #(let [link  (first (html/attr-values % :href))
5                title (html/text %)]
6            {title link})
7         (html/select (get-mainContent-html (get-html nav-link))
8                      [:div.mBd :ul.pageTPList :li :div.title :a.tit])))
9
10 (defn get-all-article-links
11   "Get a nav's all articles link and title with map as return."
```

```clojure
12    [nav-link]
13    (for [n (range 1 (inc (get-total-result-pages nav-link)))]
14      (let [url (str nav-link (format "_%d" n))]
15        (get-page-article-links url)))))
16
17 (def nav-articles-links (get-all-article-links
   ↪    "http://www.sxszjzx.com/xydt/xyxw"))
```

```clojure
1 (pprint (take 3 nav-articles-links))
```

```clojure
1 (({" 六十载同心同德建名校　一甲子匠智匠力创品牌"
   ↪    "/xydt/xyxw/content_39935"}
2   {" 你的样子，我刚好喜欢　——我校开展 "仪容仪表示范班" 评比活动"
     ↪    "/xydt/xyxw/content_39927"}
3   {" 当快闪遇上诗歌——我们一起告白我的国" "/xydt/xyxw/content_39923"}
4   {" "红色匠心，青春向党" —我校十月 "祖国颂" 诗歌朗诵比赛圆满结束"
     ↪    "/xydt/xyxw/content_39893"}
5   {" 第三届校教职工气排球圆满结束 " "/xydt/xyxw/content_39799"}
6   {" 关注课堂，携手育人——我校开展家长开放日观摩课活动报道"
     ↪    "/xydt/xyxw/content_39808"}
7   {" "更高、更快、更强" ——绍兴技师学院（筹）绍兴市职教中心成功举办第七十
     ↪    一届田径运动会"
8    "/xydt/xyxw/content_39589"}
9   {" 弘扬垦荒精神 铸牢党性之魂——我校党员教师赴大陈岛开展专题党日活动"
     ↪    "/xydt/xyxw/content_39560"}
10  {" 加强师资建设，争创名师团队——我校召开《人性的追问与教师的职业成长》
     ↪    专题讲座" "/xydt/xyxw/content_39499"}
11  {" 五星三名·用行动践行先锋力量——我校开展系列 "党员育人" 活动"
     ↪    "/xydt/xyxw/content_39157"})
```

```
12  ({" 柯桥区职教中心四位中层干部来我校交流学习"
↪    "/xydt/xyxw/content_39159"}
13   {" 江苏泰州机电高等职业技术学校一行来校学习调研"
↪     "/xydt/xyxw/content_39552"}
14   {" 倾听教师心声，助力学校发展——我校召开 2018 学年第一学期教师座谈会"
↪     "/xydt/xyxw/content_39158"}
15   {" 校工会举行新学期首次工作会议" "/xydt/xyxw/content_38999"}
16   {" "急救侠" 集结，学校安全再添 "保护伞" —我校教师积极参加 "AED" 培训"
↪     "/xydt/xyxw/content_29630"}
17   {" "年轻的战士，我为你鼓掌" ——2018 级新生军训会操纪实"
↪     "/xydt/xyxw/content_29629"}
18   {" 以奖促学——让榜样的力量发光发热" "/xydt/xyxw/content_29628"}
19   {" 未来工匠，领跑新征程——我校举行 2018 年第一学期开学典礼暨表彰大会
↪    （图文）" "/xydt/xyxw/content_29627"}
20   {" 我校学子成功晋级第三届浙江省科学玩家青少年才能挑战赛复赛"
↪     "/xydt/xyxw/content_29626"}
21   {" 告别 "压力山大"  美好从 "心" 开始—我校第一堂 "心理慕课" 顺利开课"
↪     "/xydt/xyxw/content_29625"})
22  ({" 树目标 重创新 团结协作创实效----市教育局副局长石鑫炯一行赴我校督查
↪    开学工作（图文）"
23   "/xydt/xyxw/content_29624"}
24   {" 温情期许匠心传承-我校为 2018 级新生送上最走心入学礼"
↪     "/xydt/xyxw/content_29623"}
25   {" 一个支部一个堡垒，2018 绍兴新疆班返校在路上（图文）"
↪     "/xydt/xyxw/content_29622"}
26   {" 记我校 2018 年文明单位、文明校园复评工作（图文）"
↪     "/xydt/xyxw/content_29621"}
27   {" 德育处、团委召开新学期第一次班主任会议" "/xydt/xyxw/content_34686"}
28   {" 忠于职守，继往开来——2018 学年干部集体廉政谈话（图文）"
↪     "/xydt/xyxw/content_34685"}
29  {"2018 学年编外工作人员招聘公告" "/xydt/xyxw/content_34684"}
```

```
30  {" 实践创新路上的中职名校建设——记我校 2018 学年第一期读书会（图文)"
    ↪  "/xydt/xyxw/content_34683"}
31  {" 青春服务越青商，助力拥抱大湾区" "/xydt/xyxw/content_34682"}
32  {"2018 年特招教师拟录用人员名单公示" "/xydt/xyxw/content_29640"}))
```

### 5.2.3  DONE benchmark new website two entry URL

**CLOSED:** *[2018-11-20 Tue 13:17]*

```clojure
1 (require '[clj-http.client :as http])
2 (require '[net.cgrand.enlive-html :as html])
3 (use 'criterium.core)
4
5 (pr (quick-bench (http/get "http://www.sxszjzx.com")))
```

```clojure
1 (pr (quick-bench (http/get "http://zjzx.sxsedu.net")))
```

## 5.3  TODO 从旧网站遍历所有文章，然后以标题在新网站中搜索，对比搜索结果中的第一个 [0/1]

□ 有一个统一而且简单的办法是在旧网站上遍历所有文章，然后每条文章在新网站中搜索。获得新网站的搜索结果。然后比对新旧网站文章页面的内容。

jyj.sxsedu.net/s?sid=22&wd="QUERY"

```clojure
1 ;;; Search entry
2 ;; http://jyj.sxsedu.net/s?sid=22&wd="QUERY"
3 (defn website-new-search
4   "Search TITLE in new website http://jyj.sxsedu.net/."
```

```
5    [title]
6    )
7
8    ;; div.s-result > div.result-list > ul > li > h4 > a[href="url"]
9    (defn website-root-search-result-select
10     "Select the first result from `website-new-search` results."
11     [results]
12     (-> results
13         (html/select [:.s-result :.result-list :ul :li :h4 :a])
14         (html/attr-values :href))
15     )
16
17   ;;; TODO: encode Chinese into URL
```

### 5.3.1 **TODO 如何遍历所有文章** [0/1]

☐ search

1. **DONE** 遍历爬取旧网站的所有文章 **CLOSED:** *[2018-09-27 Thu 14:33]*
   Old website 类目: http://www.sxszjzx.com/e/action/ListInfo/?classid=390
   文章页面: http://www.sxszjzx.com/e/action/ShowInfo.php?classid=390&id=
   12767

   文章页面的内容 tag

   ```
   :div.page_right > :.pright_t1 > :.pright_t4
         article         title          content
   ```

   Single entry test:

   ```
   1 (:status (http/get
   ↪    "http://www.sxszjzx.com/e/action/ShowInfo.php?classid=6&id=6"
   2                {:as "GB2312"}));; => 200
   ```

57

```clojure
3
4 (if (= " 信息提示" (first
5                    (:content
6                     (first
7                      (html/select
8                       (html/html-snippet
9                        (:body (http/get
                           ↪  "http://www.sxszjzx.com/e/action/ShowInfo.php?classid=6&id=
10                                    {:as "GB2312"})))
11                     [:title])))))
12   false
13   true)
```

A huge double deep iteration:

```clojure
1 (require 'clojure.java.io)
2
3 (for [classid (range 0 1000)]
4   (for [id (range 0 1000)]
5     (let [url     (format
       ↪  "http://www.sxszjzx.com/e/action/ShowInfo.php?classid=%s&id=%s"
6                        classid id)
7          respond (http/get url {:as "GB2312"})]
8       (if-let [status (not (= " 信息提示" (first
9                                        (:content (first
10                                          (html/select
11
                                            ↪  (html/html-snippet
                                            ↪  (:body
                                            ↪  respond))
12                                          [:title]))))))]
```

```clojure
13          (let [*out* (clojure.java.io/writer
     ↪   "/home/stardiviner/titles.lst")]
14        (println (str "OK! " url))
15        (println (-> respond
16                     html/html-snippet
17                     (html/select [:pright_t1]) ; title
18                     html/texts)))
19      (println (str "FAILED! " url))
20      ))))
```

After long time running and test, seems this does not work.

### 5.3.2   DONE 新网站搜索入口 API

**CLOSED:** *[2018-09-27 Thu 14:07]*  **DEADLINE:** *<2018-09-27 Thu>*

http://jyj.sxsedu.net/s?sid=22&wd="QUERY"

```clojure
1 (doseq [title (....)] ; TODO:
2   (-> (http/get "http://jyj.sxsedu.net/s?sid=22&wd=%s" title)
3      ))
```

### 5.3.3   DONE 获取第一个搜索结果是否和旧网站的标题相同？

**CLOSED:** *[2018-09-27 Thu 14:14]*  **DEADLINE:** *<2018-09-27 Thu>*

```
1 div.s-result > div.result-list > ul > li > h4 > a[href="url"]
```

For testing:

http://www.sxszjzx.com/e/action/ShowInfo.php?classid=115&id=15463 "致匠心"—
-艺术设计系开学第一课

http://www.sxszjzx.com/e/action/ShowInfo.php?classid=115&id=15453 17 学前教育
2 班开学第一课—君子以自强不息

http://www.sxszjzx.com/e/action/ShowInfo.php?classid=115&id=15421 童心颂师恩

```clojure
1  (require '[clojure.string :as str])
2
3  (defn website-new-search-result-same? [title]
4    "Search TITLE in new website search API."
5    (let [first-result        (first
6                                (-> (:body
7                                    (http/get
8                                     (format
                                       ↪  "http://jyj.sxsedu.net/s?sid=22&wd=%s"
                                       ↪  title)))
9                                  html/html-snippet
10                                 (html/select [:.s-result :.result-list :ul
                                     ↪  :h4 :a])))
11          first-result-title (apply str (html/texts (html/select
                ↪  first-result [:span])))
12          first-result-link  (first (html/attr-values first-result :href))]
13      (= first-result-title (convert-old-title title)))))
14
15 (defn- convert-old-title [title]
16   "Remove some special characters in old website title."
17   (str/replace title #"[ ""-]" ""))
18
19 (prn (website-new-search-result-same? " 童心颂师恩"))
20 (prn (website-new-search-result-same? " "致匠心"----艺术设计系开学第一课"))
```

```
true
true
```

## 5.4 **FAILED 按照 URL 的格式遍历所有类目和文章**

**CLOSED:** *[2018-09-25 Tue 10:45]*

Old website 类目: `http://www.sxszjzx.com/e/action/ListInfo/?classid=390` 文章页面: `http://www.sxszjzx.com/e/action/ShowInfo.php?classid=390&id=12767`

`http://www.sxszjzx.com/e/action/ListInfo/index.php?page=109&classid=33&totalnum=1645`

New website 类目: `http://zjzx.sxsedu.net/xydt/xyxw_2` 文章页面: `http://zjzx.sxsedu.net/xydt/xyxw/content_34682`

# 6 **DONE crawler.clj 爬虫存储提取连接的数据库 [5/5]**

**CLOSED:** *[2018-11-20 Tue 21:58]*

## 6.1 **DONE Redis**

**CLOSED:** *[2018-11-20 Tue 21:58]*
Clojure Redis Packages
Redis

### 6.1.1 **DONE how to access the Redis?**

**CLOSED:** *[2018-11-21 Wed 19:59]*

```clojure
1 (require '[taoensso.carmine :as redis])
2
3 ;;; Redis store crawler links
4 (defonce redis-conn-pool {:pool {}
```

```clojure
5                              :spec {:host "127.0.0.1" :port 6379}})

6

7  (defmacro wcar*

8    [& body]

9    `(redis/wcar redis-conn-pool ~@body))

10

11 (comment

12   (wcar*

13    (redis/ping)))

14

15 (defn save-link-to-redis

16   "Save crawled link to Redis.

17   Usage: (save-link-to-redis :old title link)"

18   [website-key title link]

19   (wcar*

20    (redis/hset website-key title link)))

21

22 (defn get-link-from-redis

23   "Get link to crawl from Redis.

24   Usage: (get-link-from-redis :old title"

25   [website-key title]

26   (wcar*

27    (redis/hget website-key title)))
```

### 6.1.2  DONE extract all links in nav bar

**CLOSED:** *[2018-11-22 Thu 08:44]*

```clojure
1  (defn crawl-website-old-links []

2    "Crawl all links of website old."
```

```clojure
3    (for [nav-link (map #(-> % vals first) old/nav-bar-links-map)]
4      (for [n (range 1 (inc (count old/nav-bar-links-map)))]
5        (wcar* (redis/hset :old/links n nav-link)))))

6

7 (defn crawl-website-new-links []
8   "Crawl all links of website new."
9   (for [nav-link (map #(-> % vals first) new/nav-bar-links-map)]
10     (for [n (range 1 (inc (count new/nav-bar-links-map)))]
11       (wcar* (redis/hset :new/links n nav-link)))))
```

### 6.1.3  **DONE** save all links to Redis

**CLOSED:** *[2018-11-21 Wed 21:58]*

```clojure
1 (defn crawl-website-old-links []
2   "Crawl all links of website old."
3   (for [nav-link (map #(-> % vals first) old/nav-bar-links-map)]
4     (for [n (range 1 (inc (count old/nav-bar-links-map)))]
5       (redis/hset :old/links n nav-link))))

6

7 (defn crawl-website-new-links []
8   "Crawl all links of website new."
9   (for [nav-link (map #(-> % vals first) new/nav-bar-links-map)]
10     (for [n (range 1 (inc (count new/nav-bar-links-map)))]
11       (redis/hset :new/links n nav-link))))
```

Check out the stored links:

```
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
```

```
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
```

```
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
```

```
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
```

```
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
http://www.sxti.zj.cn/html/recruit/plan.html
```

```
http://www.sxti.zj.cn/html/recruit/plan.html
```

### 6.1.4  DONE literal over all redis Hash key-values

**CLOSED:** *[2018-11-21 Wed 19:59]*

```
1 (defn all-links-in-redis
2   "Get all links in Redis."
3   [website-key]
4   (wcar*
5   ;; (redis/hgetall website-key)
6   (redis/hvals website-key)))
```

```
1 (wcar*
2  (redis/hvals :old/links))
```

# 7  TODO crawler.clj 爬取页面的内容 [3/6]

## 7.1  DONE scrape the article

**CLOSED:** *[2018-11-21 Wed 19:27]*

```
1 (defn crawl-website-old-articles []
2   "Crawl all articles of website old."
3   (for [link (all-links-in-redis :old/links)]
4     (old/get-html link)))
5
6 (defn crawl-website-new-articles []
7   "Crawl all articles of website new."
```

```
8    (for [link (all-links-in-redis :new/links)]
9      (new/get-html link)))
```

### 7.1.1  DONE old

**CLOSED:** *[2018-11-21 Wed 19:14]*

```
1  (defn get-page-article-content-html
2    "Get page's article content html."
3    [url]
4    (html/select (get-html url) [:div.page_right]))
5
6  (defn parse-article-title-and-content
7    "Parse and extract the title and content text."
8    [url]
9    (let [content-html (get-page-article-content-html url)
10         title        (html/text (first (html/select content-html
            ↪  [:div.news_Title1])))
11         content      (html/text (first (html/select content-html
            ↪  [:div.news_Content1])))]
12      {:title title, :content content}))
13
14 (comment
15   (html/text
16    (first
17     (html/select
18      (get-page-article-content
         ↪  "http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=61&id=15533")
19      [:div.news_Content1]))))
```

### 7.1.2 DONE new

**CLOSED:** *[2018-11-21 Wed 19:27]*

```clojure
1  (defn get-page-article-content-html
2    "Get page's article content html."
3    [url]
4    (html/select (get-html url) [:div.mainContent]))
5
6  (defn parse-article-title-and-content
7    "Parse and extract the title and content text."
8    [url]
9    (let [content-html (get-page-article-content-html url)
10         title        (html/text
11                        (first
12                          (html/select
13                            content-html
14                            [:div.mBd :article.articleCon :div.printArea
                                ↪  :h2.title])))
15         content      (html/text
16                        (first
17                          (html/select
18                            content-html
19                            [:div.mBd :article.articleCon :div.printArea
                                ↪  :div.conTxt])))]
20      {:title title, :content content}))
21
22  (comment
23    (html/text
24      (first
25        (html/select
```

70

```
26     (get-page-article-content-html
   ↪   "http://www.sxszjzx.com/xydt/xyxw/content_39935")
27     [:h2.title]))))
```

## 7.2  **TODO** process the crawled article content text [0/1]

### 7.2.1  **TODO** strip some special characters and spaces

☐

☐ _

## 7.3  **TODO** auto find content element

How to auto find the max scope content element's text. So that crawler can auto extract the text. No need the selector rules.

# 8  **TODO** store.clj 爬虫内容存储到数据库 [3/4]

## 8.1  **DONE** SQLite

**CLOSED:** *[2018-11-22 Thu 10:28]*

```clojure
1 (defn save-to-sqlite
2   "Save data into SQLite DB."
3   [title article]
4   (jdbc/with-db-connection [db sqlite]
5     (jdbc/insert! db :articles {:title   title ; TODO: do I need insert
   ↪   the `:id`?
6                                 :article article})))
7
8
9 (defn save-all-articles-to-sqlite
```

71

```
10    []
11    (for [article (crawler/crawl-website-old-articles)]
12      (save-to-sqlite (key article) (val article))))
```

### 8.1.1  DONE detect table exist?

**CLOSED:** *[2018-11-22 Thu 10:08]*

Use {:conditional? true} for jdbc/create-table-ddl.

### 8.1.2  DONE save title and content to SQLite

**CLOSED:** *[2018-11-22 Thu 10:28]*

```
1  (defn all-links-in-redis
2    "Get all links in Redis."
3    [website-key]
4    (wcar*
5     ;; (redis/hgetall website-key)
6     (redis/hvals website-key)))
7
8  (defn crawl-website-old-articles []
9    "Crawl all articles of website old."
10   (crawl-website-old-articles)
11   (for [link (all-links-in-redis :old/links)]
12     (let [page (old/parse-article-title-and-content link)]
13       (store/save-to-sqlite :old (:title page) (:content page)))))
14
15 (comment
16   (let [kk (old/parse-article-title-and-content (first (all-links-in-redis
     ↪  :old/links)))]
17     (store/save-to-sqlite :old (:title kk) (:content kk))))
```

```
18
19 (defn crawl-website-new-articles []
20   "Crawl all articles of website new."
21   (crawl-website-new-links)
22   (for [link (all-links-in-redis :new/links)]
23     (let [page (new/parse-article-title-and-content link)]
24       (store/save-to-sqlite :new (:title page) (:content page)))))
25
26 (comment
27   (let [kk (new/parse-article-title-and-content (first (all-links-in-redis
     ↪  :new/links)))]
28     (store/save-to-sqlite :new (:title kk) (:content kk))))
```

## 8.2 **TODO** MongoDB

# 9 **TODO** validate.clj (检查页面的可用性)

## 9.1 Validate Links

## 9.2 Validate Images

## 9.3 Validate Videos

## 9.4 Validate Docs

# 10 **TODO** compare.clj 对比文章

☒ 对比文章标题

☐ 对比 html 页面中的 img，video 之类的 tags 的数量和类型。

## 10.1 **TODO** compare.clj 比较页面的相似性 [0/1]

☐ 对比文字部分 tag 的内容，用 MD5 之类的 digest 算法

- SimHash

- MinHash

☐ search Simhash Java libraries Java MinHash

# 11  DONE 比较上传视频

**CLOSED:** *[2018-09-26 Wed 19:08]*

## 11.1  DONE 旧网站的视频

**CLOSED:** *[2018-09-26 Wed 19:08]*

```clojure
1 (def website-old-url "http://www.sxti.zj.cn")
```

nil#'user/website-new-url

get a video link in list example:

```clojure
1 (pprint
2  (first
3   (:content
4    (first
5     (html/select
6      (-> (:body
7           (http/get
8            (format (str website-old-url
                 ↪ "/e/action/ListInfo/index.php?classid=10&page=%s") 0)
9            {:as "GB2312"}))
10          html/html-snippet
11          (html/select [:.imgList1]))
12      [:ul]))))))
```

74

```
1  {:tag :li,
2   :attrs nil,
3   :content
4   ({:tag :a,
5     :attrs
6     {:href "/e/action/ShowInfo.php?classid=173&id=15401",
7      :target "_blank"},
8     :content
9     ({:tag :img,
10       :attrs
11       {:src

12
       ↪   "/d/file/html/news/4/4/2018-09-07/small29350dd5474f452348699de611ad62f8.jpg",
13       :alt " 心理课堂六",
14       :width "220",
15       :height "180"},
16       :content nil})}
17   {:tag :div, :attrs {:class "movie1"}, :content nil}
18   {:tag :h1,
19     :attrs nil,
20     :content
21     ({:tag :a,
22       :attrs
23       {:href "/e/action/ShowInfo.php?classid=173&id=15401",
24        :title " 心理课堂六",
25        :target "_blank"},
26       :content (" 心理课堂六")})})})}
```

```
1 {:tag      :li,
2  :attrs    nil,
3  :content ({:tag       :a,
4             :attrs    {:href "/e/action/ShowInfo.php?classid=378&id=12761",
              ↪  :target "_blank"},
5             :content ({:tag :img, :attrs {:src
              ↪  "/d/file/html/news/4/5/2017-05-24/small9bfeabf676860034d33a510fcd583c35.jpg
              ↪  :alt " 重拾自我---16 旅服高工", :width "220", :height
              ↪  "180"}, :content nil})
6             }
7            {:tag       :div,
8             :attrs    {:class "movie1"},
9             :content nil}
10           {:tag       :h1,
11            :attrs    nil,
12            :content ({:tag       :a,
13                       :attrs    {:href
                         ↪  "/e/action/ShowInfo.php?classid=378&id=12761",
                         ↪  :title " 重拾自我---16 旅服高工", :target
                         ↪  "_blank"}
14                       :content (" 重拾自我---16 旅服高工")})
15           })}
```

```
1 ;;; get first list page for testing
2 (comment
3   (def page-list
4     (-> (:body
5         (http/get
6          (format (str website-old-url
           ↪  "/e/action/ListInfo/index.php?classid=10&page=%s") 0)
```

```clojure
7         {:as "GB2312"}))
8       html/html-snippet
9       (html/select [:.imgList1]))))

11 (defn- extract-video-page-content [page]
12   (-> (:body
13       (http/get
14         (format (str website-old-url
              ↪   "/e/action/ListInfo/index.php?classid=10&page=%s") page)
15         {:as "GB2312"}))
16     html/html-snippet
17     (html/select [:.imgList1])))

19 ;;; get all title and URL
20 (defn- extract-video-link [node]
21   (let [lists ((juxt #(first (html/attr-values % :href))
22                      #(first (html/attr-values % :title)))
23               (-> node
24                   (html/select [:h1 :a])
25                   first))]
26     {(second lists)
27      ;; construct url
28     (str website-old-url (first lists))}))

30 (comment
31   (map
32    extract-video-link
33    (html/select-nodes* (html/select page-list [:ul]) [:li])))
```

```
1  (defn extract-page-videos-maxpage-old
2    "Get the max page number of
       ↪ http://www.sxti.zj.cn/e/action/ListInfo/?classid=10"
3    [url]
4    (let [last_page (-> (:body (http/get url {:as "GB2312"}))
5                        html/html-snippet
6                        (html/select [:div.yema1])
7                        first
8                        :content
9                        last)]
10     (if (= (first (:content last_page)) " 尾页")
11       ;; (print (format "http://www.sxti.zj.cn/%s" (first
              (html/attr-values last_page :href))))
12       (Integer. (first (re-seq #"\d+"
13                                (re-find #"page=\d+&" (first
                                   ↪ (html/attr-values last_page
                                   ↪ :href)))))))
14     ))
15
16 (print (extract-page-videos-maxpage-old (str website-old-url
   ↪ "/e/action/ListInfo/?classid=10")))
```

13

#+RESULTS[*<2018-10-16 Tue 14:51>* bdecab31eb8bad0dfc2c2d281094d39fe174346c]:
extract-video-pages-max-old

```
1 <<define-old-website-url>>
2 <<extract-video-page-info-old>>
3 <<extract-video-pages-max-old>>
4
```

```
5  (def videos-links-old
6    (reduce
7     merge
8     (reduce
9      concat
10     (for [n (range 0 (extract-page-videos-maxpage-old (str website-old-url
       ↪  "/e/action/ListInfo/?classid=10")))]
11       (map
12        #(-> %
13             extract-video-link)
14        (html/select-nodes*
15         (html/select (extract-video-page-content n) [:ul])
16         [:li]))))))))
17
18 (pprint videos-links-old)
```

## 11.2  DONE 新网站的视频

**CLOSED:** *[2018-09-26 Wed 19:08]*

```
1  (def website-new-url "http://www.sxszjzx.com/")
2  ;; (def website-new-url "http://zjzx.sxsedu.net")
```

nil#'user/website-new-url

get a video link in list example:

```
1  (pprint
2    (first
3     (html/select
```

```clojure
   (-> (:body
        (http/get
         (format (str website-new-url "/ztlm/xysp_%s") 1)))
       html/html-snippet
       (html/select [:ul.vedioPageList]))
   [:li :h3 :a])))
```

```clojure
{:tag :a,
 :attrs {:href "/ztlm/xysp/zycz/content_129", :title "《青春 梦想》 14 大专
↪   动漫"},
 :content ("《青春 梦想》 14 大专动漫")}
```

```clojure
(defn extract-page-videos-maxpage-new
  "Get the max page number of http://www.sxszjzx.com/ztlm/xysp."
  [url]
  (Integer.
   (clojure.string/replace
    (re-find #"/*\d+ 页"
             (-> (:body (http/get url))
                 html/html-snippet
                 (html/select [:div.page :span.total])
                 first
                 :content
                 first))
    " 页" "")))

(print (extract-page-videos-maxpage-new (str website-new-url
↪   "/ztlm/xysp")))
```

14

#+RESULTS[*<2018-10-16 Tue 14:53>* 92e10ba363a1a208d764be2f96efcc2f1e4a5be4]: extract-video-pages-max-new

```
1  <<define-new-website-url>>
2  <<extract-video-pages-max-new>>
3
4  (def videos-links-new
5    (reduce
6     merge
7     (reduce
8      concat
9      (for [n (range 1 (extract-page-videos-maxpage-new (str website-new-url
       ↪  "/ztlm/xysp")))]
10       (map
11        (fn [a]
12          {(first (html/attr-values a :title))
13           (str "http://zjzx.sxsedu.net" (first (html/attr-values a
             ↪  :href)))})
14        (html/select
15         (-> (:body
16              (http/get
17               (format "http://zjzx.sxsedu.net/ztlm/xysp_%s" n)))
18             html/html-snippet
19             (html/select [:ul.vedioPageList]))
20         [:li :h3 :a]))))))
21
22 (pprint videos-links-new)
```

14{"纪念刘和珍君" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_106",
 "走近绍兴黄酒1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_76",
 "《青春" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_129",

"人生因设计而美丽2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_90",
"3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_94",
"基本站姿组合手位1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_84",
"计数原理" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_97",
"《展望未来我们的明天会更加精彩》14大专城建"
"http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_132",
"Art" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_93",
"动画动作修饰2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_86",
"绞孔攻螺纹3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_83",
"方言与普通话2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_107",
"07-课堂礼仪" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_119",
"色彩空间" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_101",
"局域网配置综合实训3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_75",
"——有一种爱叫做亲情" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_37",
"我被十三所学校开除(2)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_28",
"中央电视台-开学第一课" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_141",
"人生因设计而美丽1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_89",
"王金云" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_128",
"前厅模拟接待" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_66",
"人生因设计而美丽3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_91",
"重拾自我---16旅服高工" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_136",
"Body" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_63",
"2017年学第二学期开学第一课" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_61",
"The" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_104",
"《拾梦者》" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_138",
"室内操" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_21",
"懂得体贴孝敬父母 - 下" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_112",
"第二届三字经PK赛" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_41",
"08-校园礼仪" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_120",
"12-社会交往（下）" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_124",
"绞孔攻螺纹1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_81",

"我喜欢我的学校（校园生活纪录）" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_51",

"葛豪焙" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_135",

"孝亲尊师电视特别节目（三）" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_39",

"礼仪操" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_46",

"05-教师的装饰" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_117",

"永恒的舞台" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_53",

"全国人口普查中小学生一堂课动画短剧"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_31",

"纺织贸易实务" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_68",

"心理课堂六" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_126",

"我和我的学校（浙江省影视大赛第一名作品）"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_52",

"社团成果展示片" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_49",

"2014预科生活——预习青春" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_59",

"01-教师礼仪概述" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_113",

"孝亲尊师电视特别节目(二)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_38",

"绍兴市职教中心建校50周年校庆" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_139",

"第三套中学生广播体操-舞动青春分解示范"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_23",

"商务学区班主任论坛2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_96",

"一个人---15大专动漫" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_137",

"合情推理3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_80",

"商务学区班主任论坛1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_95",

"直系电机系统故障排除" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_69",

"礼仪教育视频" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_125",

"健康消费\"茶文化知识讲座(2)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_33",

"开学第一课（2011年9月1日）"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_140",

"14计算机1班" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_134",

"健康消费\"茶文化知识讲座(3)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_34",

"《青春誓言》13美术高考2班" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_133",

"绞孔攻螺纹2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_82",
"《技能成就梦想》13大专机电徐俊亮、傅卓楷"
"http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_131",
"第五十八届田径运动会" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_43",
"职教三字经PK赛1" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_26",
"10-师生关系" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_122",
"2017新生军训纪录片" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_60",
"健康消费\"茶文化知识讲座(4)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_35",
"《正步人生》绍兴市职教中心2010军训特别节目"
"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_142",
"邹越：让生命充满爱" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_29",
"基本站姿组合手位3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_88",
"懂得体贴孝敬父母 - 中" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_111",
"合情推理1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_72",
"绍兴市职教中心2011军训特别节目"
"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_40",
"艺术教育特色宣传片" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_25",
"04-教师的语言" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_116",
" "国家改革发展示范校"建设纪实" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_55",
"徐涵宗" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_130",
"\"走进\"世贸" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_100",
"房志成" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_127",
"11-社会交往（上）" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_123",
"我被十三所学校开除(1)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_27",
"方言与普通话3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_70",
"茶为国饮" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_36",
"职教之歌" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_56",
"学会节俭，不做奢侈的孩子 - 上" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_108",
"绍兴市职教中心" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_18",
"第三套中学生广播体操-舞动青春完整音乐"
"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_24",

"方言与普通话1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_105",

"06-教师的仪表" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_118",

"局域网配置综合实训2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_74",

"合情推理2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_78",

"走近绍兴黄酒2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_77",

"学会节俭，不做奢侈的孩子 - 下" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_109",

"神神慢" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_103",

"走近绍兴黄酒3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_79",

"2007校园文化节" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_19",

"09-办公室礼仪" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_121",

"石膏像的打形" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_67",

"新兴力量的崛起" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_98",

"多彩的消费" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_64",

...}

## 11.3 DONE 两个视频数据集合比较

**CLOSED:** *[2018-09-26 Wed 19:08]*

```
1 (require '[clojure.set :as set])
2
3 <<website-old-videos>>
4 <<website-new-videos>>
5
6 (println "----------------------------------")
7
8 (map
9  println
10  (set/difference
11   (set (keys videos-links-old))
12   (set (keys videos-links-new)))))
```

13{nil "http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5601",

"十月主题演讲比赛"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=13200",

"纪念刘和珍君"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=6427",

"走近绍兴黄酒1"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5658",

"《青春"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11369",

"人生因设计而美丽2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5638",

"3" "http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5687",

"基本站姿组合手位1"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5634",

"计数原理"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5599",

"《展望未来我们的明天会更加精彩》14大专城建"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11373",

"Art"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5685",

"动画动作修饰2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5665",

"绞孔攻螺纹3"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5663",

"方言与普通话2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5669",

"07-课堂礼仪"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5531",

"色彩空间"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5603",

"局域网配置综合实训3"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5655",
"中央电视台-开学第一课"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5564",
"人生因设计而美丽1"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5637",
"王金云"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11368",
"我的学校（专业选报指导）"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=6434",
"人生因设计而美丽3"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5639",
"重拾自我---16旅服高工"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=12761",
"2017年学第二学期开学第一课"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=14247",
"The"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5598",
"《拾梦者》"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11375",
"室内操"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5628",
"懂得体贴孝敬父母 - 下"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5572",
"08-校园礼仪"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5532",
"预防电信网络诈骗专题讲座"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=13151",
"12-社会交往（下）"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5536",
"绞孔攻螺纹1"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5661",
"我喜欢我的学校（校园生活纪录）"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=6438",
"葛豪焙"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11376",
"我被十三所学校开除1"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5558",
"母亲节"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=9315",
"礼仪操"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=6428",
"05-教师的装饰"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5529",
"永恒的舞台"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=6440",
"全国人口普查中小学生一堂课动画短剧"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5556",
"纺织贸易实务"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5751",
"心理课堂六"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=15401",
"我和我的学校（浙江省影视大赛第一名作品）"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=6439",
"社团成果展示片"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=6436",
"01-教师礼仪概述"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5521",
"孝亲尊师电视特别节目(二)"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5496",
"绍兴市职教中心建校50周年校庆"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5630",

"第三套中学生广播体操-舞动青春分解示范"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5612",
"共创美好未来"大合唱比赛决赛"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=15167",
"商务学区班主任论坛2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5641",
"一个人---15大专动漫"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=12762",
"合情推理3"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5633",
"商务学区班主任论坛1"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5640",
"直系电机系统故障排除"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5755",
"礼仪教育视频"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=6429",
"开学第一课（2011年9月1日）"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5491",
"14计算机1班"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11374",
"2017校园新歌声总决赛"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=13711",
"《青春誓言》13美术高考2班"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11372",
"绞孔攻螺纹2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5662",
"《技能成就梦想》13大专机电徐俊亮、傅卓楷"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11371",
"第五十八届田径运动会"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5831",
"职教三字经PK赛1"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5589",
"10-师生关系"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5534",
"2017新生军训纪录片"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=13117",
"《正步人生》绍兴市职教中心2010军训特别节目"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5557",
"邹越：让生命充满爱"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5560",
"2014预科生活 预习青春"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=9374",
"基本站姿组合手位3"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5636",
"懂得体贴孝敬父母 - 中"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5571",
"我被十三所学校开除2"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5559",
"合情推理1"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5631",
"绍兴市职教中心2011军训特别节目"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5440",
"艺术教育特色宣传片"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5607",
"04-教师的语言"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5528",
""国家改革发展示范校"建设纪实"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=8017",
"徐涵宗"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11370",
"房志成"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=378&id=11367",

"11-社会交往（上）"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5535",
"方言与普通话3"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5670",
"我的学校 绍兴市职教中心"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=6435",
"茶为国饮"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5550",
"职教之歌"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=8288",
"学会节俭，不做奢侈的孩子 - 上"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5566",
"第三套中学生广播体操-舞动青春完整音乐"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5613",
"方言与普通话1"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5668",
"06-教师的仪表"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5530",
"职教三字经PK赛2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5590",
"局域网配置综合实训2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5654",
"合情推理2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5632",
"走近绍兴黄酒2"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5659",
"学会节俭，不做奢侈的孩子 - 下"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5569",
"神神慢"

"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5604",
"走近绍兴黄酒3"

```
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5660",
"09-办公室礼仪"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=173&id=5533",
"新兴力量的崛起"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5600",
"藏族民族风情"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5605",
"07级新生军训"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5828",
"动画动作修饰3"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=172&id=5666",
"孝亲尊师电视特别节目(二)成长足迹"
"http://www.sxti.zj.cn/e/action/ShowInfo.php?classid=171&id=5495",
...}
14{"纪念刘和珍君" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_106",
"走近绍兴黄酒1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_76",
"《青春" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_129",
"人生因设计而美丽2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_90",
"3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_94",
"基本站姿组合手位1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_84",
"计数原理" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_97",
"《展望未来我们的明天会更加精彩》14大专城建"
"http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_132",
"Art" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_93",
"动画动作修饰2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_86",
"绞孔攻螺纹3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_83",
"方言与普通话2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_107",
"07-课堂礼仪" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_119",
"色彩空间" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_101",
"局域网配置综合实训3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_75",
"——有一种爱叫做亲情" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_37",
```

"我被十三所学校开除(2)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_28",
"中央电视台-开学第一课" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_141",
"人生因设计而美丽1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_89",
"王金云" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_128",
"前厅模拟接待" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_66",
"人生因设计而美丽3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_91",
"重拾自我---16旅服高工" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_136",
"Body" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_63",
"2017年学第二学期开学第一课" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_61",
"The" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_104",
"《拾梦者》" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_138",
"室内操" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_21",
"懂得体贴孝敬父母 - 下" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_112",
"第二届三字经PK赛" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_41",
"08-校园礼仪" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_120",
"12-社会交往（下）" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_124",
"绞孔攻螺纹1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_81",
"我喜欢我的学校（校园生活纪录）" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_51",
"葛豪焙" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_135",
"孝亲尊师电视特别节目（三）" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_39",
"礼仪操" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_46",
"05-教师的装饰" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_117",
"永恒的舞台" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_53",
"全国人口普查中小学生一堂课动画短剧"
"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_31",
"纺织贸易实务" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_68",
"心理课堂六" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_126",
"我和我的学校（浙江省影视大赛第一名作品）"
"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_52",
"社团成果展示片" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_49",
"2014预科生活——预习青春" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_59",

"01-教师礼仪概述" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_113",

"孝亲尊师电视特别节目(二)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_38",

"绍兴市职教中心建校50周年校庆" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_139",

"第三套中学生广播体操-舞动青春分解示范"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_23",

"商务学区班主任论坛2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_96",

"一个人---15大专动漫" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_137",

"合情推理3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_80",

"商务学区班主任论坛1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_95",

"直系电机系统故障排除" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_69",

"礼仪教育视频" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_125",

"健康消费\"茶文化知识讲座(2)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_33",

"开学第一课（2011年9月1日）"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_140",

"14计算机1班" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_134",

"健康消费\"茶文化知识讲座(3)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_34",

"《青春誓言》13美术高考2班" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_133",

"绞孔攻螺纹2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_82",

"《技能成就梦想》13大专机电徐俊亮、傅卓楷"

"http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_131",

"第五十八届田径运动会" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_43",

"职教三字经PK赛1" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_26",

"10-师生关系" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_122",

"2017新生军训纪录片" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_60",

"健康消费\"茶文化知识讲座(4)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_35",

"《正步人生》绍兴市职教中心2010军训特别节目"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_142",

"邹越：让生命充满爱" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_29",

"基本站姿组合手位3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_88",

"懂得体贴孝敬父母 - 中" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_111",

"合情推理1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_72",

"绍兴市职教中心2011军训特别节目"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_40",

"艺术教育特色宣传片" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_25",

"04-教师的语言" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_116",

" "国家改革发展示范校" 建设纪实" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_55",

"徐涵宗" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_130",

"\"走进\"世贸" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_100",

"房志成" "http://zjzx.sxsedu.net/ztlm/xysp/zycz/content_127",

"11-社会交往（上）" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_123",

"我被十三所学校开除(1)" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_27",

"方言与普通话3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_70",

"茶为国饮" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_36",

"职教之歌" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_56",

"学会节俭，不做奢侈的孩子 - 上" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_108",

"绍兴市职教中心" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_18",

"第三套中学生广播体操-舞动青春完整音乐"

"http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_24",

"方言与普通话1" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_105",

"06-教师的仪表" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_118",

"局域网配置综合实训2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_74",

"合情推理2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_78",

"走近绍兴黄酒2" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_77",

"学会节俭，不做奢侈的孩子 - 下" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_109",

"神神慢" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_103",

"走近绍兴黄酒3" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_79",

"2007校园文化节" "http://zjzx.sxsedu.net/ztlm/xysp/xnhd/content_19",

"09-办公室礼仪" "http://zjzx.sxsedu.net/ztlm/xysp/jysp/content_121",

"石膏像的打形" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_67",

"新兴力量的崛起" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_98",

"多彩的消费" "http://zjzx.sxsedu.net/ztlm/xysp/jxfc/content_64",

...}

---------------------------------------

```
nil
```

十月主题演讲比赛

我的学校（专业选报指导）

预防电信网络诈骗专题讲座

我被十三所学校开除1

母亲节

共创美好未来”大合唱比赛决赛

2017校园新歌声总决赛

2014预科生活 预习青春

我被十三所学校开除2

我的学校 绍兴市职教中心

职教三字经PK赛2

孝亲尊师电视特别节目(二)成长足迹

孝亲尊师电视特别节目(一)---有一种爱叫做亲情

## 11.4  DONE 检查新网站视频的缩略图是否存在

**CLOSED:** *[2018-10-11 Thu 15:00]*  **SCHEDULED:** *<2018-10-11 Thu>*

failed thumbnail example: `http://www.sxszjzx.com/ztlm/xysp_4`

$\\$

($_{\text{Common}}\backslash\backslash\backslash$-translation$_{240160}$.gif

---

```html
1 <li>
2     <div class="pic">
3         <a href="/ztlm/xysp/xnhd/content_43">
4             <img alt=" 第五十八届田径运动会"
              ↪  src="\Content\_Common\Base\img\error-translation_240_160.gif"
              ↪  width="240" height="160">
5         </a>
6     </div>
7     <h3 class="title">
```

```
8        <a href="/ztlm/xysp/xnhd/content_43" title=" 第五十八届田径运动
    ↪   会">第五十八届田径运动会</a>
9    </h3>
10   <span class="titleBg"></span>
11 </li>
```

---

```clojure
1 (defn extract-page-video-links
2   "Extract all video links from web page."
3   [url]
4   ;; div>a
5   (html/select
6    (-> (:body
7         (http/get url))
8       html/html-snippet
9       (html/select [:ul.vedioPageList]))
10   [:li :div.pic :a]))
11
12 (defn extract-page-video-filename
13   "Extract video filename from web page."
14   [url]
15   (last
16    (clojure.string/split
17     (first
18      (html/attr-values
19       (first
20        (-> (:body (http/get url))
21            html/html-snippet
22            (html/select [:div.vedioPlayer :div])))
23       :data-url))
24     #"video/")))
```

```clojure
25

26 (defn extract-page-videos-maxpage
27   "Get the max page number of http://www.sxszjzx.com/ztlm/xysp."
28   [url]
29   (Integer.
30    (clojure.string/replace
31     (re-find #"/*\d+ 页"
32               (-> (:body (http/get url))
33                   html/html-snippet
34                   (html/select [:div.page :span.total])
35                   first
36                   :content
37                   first))
38     " 页" "")))

39

40 ;; (extract-page-videos-maxpage "http://www.sxszjzx.com/ztlm/xysp")
```

```clojure
1 <<extract page videos>>

2

3 (for [n (range 1 (extract-page-videos-maxpage-new
  ↪  "http://www.sxszjzx.com/ztlm/xysp"))]
4   (map
5    (fn [a]
6      (let [img           (first (html/select a [:img]))
7            img-src       (first (html/attr-values img :src))
8            title         (first (html/attr-values img :alt))
9            video-page-url (str "http://www.sxszjzx.com" (first
              ↪  (html/attr-values a :href)))]
10        (when (and (not (clojure.string/blank? img-src))
11                   (re-find #"\\Content*" img-src))
```

```
12          ;; get video info link
13          (print {title video-page-url})
14          ;; get video filename
15          (println
16           (extract-page-video-filename
17            (str "http://www.sxszjzx.com" (first (html/attr-values a
                ↪   :href)))))))))))
18    (extract-page-video-links (format "http://www.sxszjzx.com/ztlm/xysp_%s"
      ↪   n))))
```

#+RESULTS[*<2018-10-16 Tue 15:00>* 4c0092ec2c05b70d98b0ba30aae6b8f2ce3c5347]:

## 11.5   DONE 获取新网站视频页面的视频源信息

**CLOSED:** *[2018-10-21 Sun 11:02]*

http://www.sxszjzx.com/ztlm/xysp/zycz/content_129

```
1  <<define-new-website-url>>
2
3  (defn website-new-video-page-extract-info
4    [url]
5    (let [video-object (-> (:body (http/get url))
6                        html/html-snippet
7                        (html/select [:div.vedioPlayer :div])
8                        first)]
9      {(apply str (html/attr-values video-object :data-title))
10      (str website-new-url (first (html/attr-values video-object
         ↪   :data-url)))}))
11
12  (website-new-video-page-extract-info
    ↪   "http://www.sxszjzx.com/ztlm/xysp/zycz/content_129")
```

```
{"《青春梦想》14大专动漫"
 "http://www.sxszjzx.com//upload/sxszjzx/contentmanage/video/2018/09/27/qcmx[自定义].mp4"}
```

## 11.6   **TODO 检查视频是否可以播放** [4/5]

☒ 检测视频源文件是否存在?

    ☒ 像 这个视频连接源文件地址 就是可以打开播放的, 但是用新网站的播放器就是不行。

    ☒ 应该是需要 flash player, 但是我检测我安装了 flash plugin 了。

GET http://www.sxszjzx.com/Content/_Common/Assets/Scripts/swfobject.js?
_=1540091532521 GET http://www.sxszjzx.com/Content/_Common/Base/swf/PowerPlayback.
swf

When I click play button on video frame:

There is a request return 404 result status:

```
http://www.sxszjzx.com/upload/sxszjzx/contentmanage/video/2018/09/27/qcmx[
```

☐ 这里应该给他们报告个 Bug, 在 Windows 下请求的 URL 是正确的, 在 Linux 下却是缺失部分文件名的, 看来是无法正确处理一些文件名中的特殊字符。

☒ 在 Linux Firefox/Firefox Developer Edition 两个浏览器下测试, 因为其他的浏览器都不支持 Flash。

request headers:

```
1 Host: www.sxszjzx.com
2 User-Agent: Mozilla/5.0 (X11; Linux x86_64; rv:62.0) Gecko/20100101 Firefox/62.0
3 Accept: text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8
4 Accept-Language: en-US,en;q=0.7,zh-CN;q=0.3
5 Accept-Encoding: gzip, deflate
6 Referer: http://www.sxszjzx.com/Content/_Common/Base/swf/PowerPlayback.swf
7 Cookie: PowerUniqueVisitor=c3f5685d-dfca-46fb-86d9-b58efc5eef6b_2018%2F10%2F21%200%3A00%3A0
```

```
8 DNT: 1
9 Connection: keep-alive
```

response headers:

```
1 HTTP/1.1 404 Not Found
2 Content-Type: text/html
3 ServerResponseDuration: 15.6279ms
4 X-PowerEasy-Version: 1.9.1.2
5 X-PowerEasy-Product: SiteAzure
6 X-Frame-Options: SAMEORIGIN
7 X-XSS-Protection: 1; mode=block
8 X-Content-Type-Options: nosniff
9 Date: Sun, 21 Oct 2018 03:15:39 GMT
10 Content-Length: 1163
```

## Usage
FIXME
## License
Copyright © 2018 FIXME

Distributed under the Eclipse Public License either version 1.0 or (at your option) any later version.

# 12   DONE Use Git as Version Control System

**CLOSED:** *[2018-11-23 Fri 08:28]*

## 13  Generate Project Report

### 13.1  update project clocktable dynamic block

### 13.2  **TODO** Org Mode export