AUTOCRAWLER : A Progressive Understanding Web Agent for Web Crawler Generation

Anonymous ACL submission

Abstract

001 Web automation is a significant technique that accomplishes complicated web tasks through automating common web actions, which enhance operational efficiency and reduce the need for manual intervention. Traditional methods, such as wrappers, suffer from limited adaptability and scalability when facing new website. On the other hand, generative agents empowered by Large Language Models (LLMs) exhibit poor performance and reusability in open-world scenarios. In this work, we introduce a crawler generation task 013 for vertical information web pages and the paradigm of combining LLMs with crawlers, which helps crawlers handle diverse and changing web environments more efficiently. We propose AUTOCRAWLER, a two-stage frame-017 work that leverages the hierarchical structure of HTML for progressive understanding. Through top-down and step-back operations, AUTOCRAWLER can learn from erroneous ac-021 tions and continuously prune HTML for better action generation. We conduct comprehensive experiments with multiple LLMs and demonstrate the effectiveness of our framework.

1 Introduction

026

Web automation refers to the process of programmatically interacting with web-based applications or websites to execute tasks that would typically require human intervention. Web automation streamlines repetitive and time-consuming tasks, significantly enhancing efficiency, accuracy, and scalability across diverse online processes.

In traditional web automation, methods predominantly rely on wrappers, which are scripts or software specifically designed to extract data from predetermined websites or pages. This approach is characteristic of a closed-world scenario, where the automation system only interacts with a predefined, limited set of websites or pages and does not extend beyond this specified domain. Consequently, these traditional methods face limitations



Figure 1: **Top:** HTML have a hierarchical structure DOM tree; **Down:** Existing web automation framework: Green arrows refer to handcraft/LLMs prompting process, Violet arrows refer to parser.

in *adaptability* and *scalability*, struggling to function effectively when encountering new or altered website structures. Given these limitations, both rule-based wrappers and auto wrappers (Bronzi et al., 2013), despite their differences, share a common dependency on manually-annotated examples for each website (Gulhane et al., 2011).

The advent of Large language models (LLMs) has revolutionized web automation by introducing advanced capabilities such as planning, reasoning, reflection and tool using. Leveraging these capabilities, web automation employs LLMs to construct generative agents that can autonomously navigate, interpret, and interact with web content, This effectively solves open-world web-based tasks through sophisticated language understanding and decisionmaking processes. However, despite these advancements, this paradigm faces two major issues. On one hand, existing web agent frameworks often demonstrate poor performance, with a success rate mentioned as 2.0 (Deng et al., 2023) on open-world tasks. On the other hand, a significant shortcoming encountered in this approach is its insufficient reusability. This implies that these agents are overly dependent on LLMs even when dealing with similar tasks, leading to low efficiency when managing a large volume of repetitive and similar webpages.

059

060

063

064

067

072

077

084

100

102

103

104

105

106

107

108

In this work, we propose a crawler generation task for vertical information web pages. The goal of this task is to automatically generate a series of predefined rules or action sequence to automatically extract target information. This task calls for a paradigm that combines LLMs and crawlers. Compared to traditional wrappers, this paradigm can be quickly adjusted according to different websites and task requirements. This flexibility enables crawlers to handle diverse and changing web environments more efficiently. Compared to the generative agent paradigm, it introduces intermediate rules to enhance reusability and reduce the dependency on LLMs when dealing with similar tasks, thereby improving efficiency when handling a large number of web tasks.

Although LLMs possess strong web page comprehension abilities, they cannot reliably understand and parse web page structures for three main reasons. First, LLMs are primarily pre-trained on massive corpora of cleansed, high-quality pure text, lacking exposure to markup languages such as HTML. As a result, LLMs exhibit a limited understanding of the complex structures and semantics inherent in HTML. Second, HTML, as a semi-structured data, amalgamates elements of both structured (tags and attributes) and unstructured (textual content), concurrently encompassing multilayered information nested. This amalgamation augment the complexity of understanding. Third, although LLMs excel in comprehending textual content, they still fall short in understanding and maintaining the structural information of lengthy documents. This indicates a potential challenge in accurately capturing and utilizing the hierarchical structure inherent in long HTML documents.

Therefore, we introduce AUTOCRAWLER, a two-phase framework designed to address the

crawler generation task. An overview of AU-109 TOCRAWLER is presented in Figure 2. Our frame-110 work leverages the hierarchical structure of HTML 111 for progressive understanding. Specifically, we pro-112 pose a heuristic algorithm based on LLMs, which 113 utilizes the DOM-tree hierarchical structure of web 114 pages, to automatically correct erroneous actions 115 and continuously optimize by pruning the irrelevant 116 parts of the HTML content. 117

Our contributions can be summarized as follow:

• We propose the web crawler generation task and the paradigm of combining LLMs and crawler together. 118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

- We introduce AUTOCRAWLER, a two-phase framework with progressive understanding to generate executable action sequences.
- Comprehensive experimental results demonstrate the effectiveness of our framework in the web crawler generation task.

2 Preliminaries

In this section, we first define the crawler generation task, and then present the dataset collection process and its corresponding evaluation metrics.

2.1 Task Formulation

We first formulate our crawler generation task. Given a set of webpages from the same website $w \in W$ describing a subject entity *s* (also called topic entity in the previous literature), and its corresponding predefined target attribute $r \in \mathcal{R}$, the goat of the task is to generate an executable rule/action sequence \mathcal{A} to extract the target information *o* from every webpages.

2.2 Datasets

We adopt semi-structure information extraction task as the testbed for crawler generation task.

SWDE DATASET (Hao et al., 2011) is an Structured Web Data Extraction dataset that contains webpages and golden label from 80 websites in 8 domains, with 124,291 webpages. Each of the websites from the same domains focus on 3-5 attributes in the webpages. Detail information is in Appendix A.

We transform the dataset with follow settings. First, we design instruction for each of the domains,



Figure 2: Our framework for generation of crawler. Left: Progressive generation process, consist of a cycle of top-down and step-back to progressively generate an executable action sequence; **Right:** Synthesis process, generate a stable action sequence generated based on ones from seed websites.

and for each of the attributes as the input information for LLMs¹. Second, for each websites in each domains, we sample 100 webpages as the whole test set. We regard the set of the webpages in the same websites and the corresponding extraction instruction as a case. For example, for the ESPN websites² in NBA player domains, the sampled 100 detail webpage of players and the instruction *Please extract the team of the player he play now* is a complete test case of our crawler generation task. Third, we preprocess the websites by removing those elements in a webpage that do not contribute to the semantics. We filter out all DOM element nodes with <script> and <style>, as well as delete all attributes in element node except @class. And we replace the original escape characters in the annotations to ensure consistency with the corresponding information on the web.

Ultimately, we collect a dataset containing 320 test cases, covering 32 different extraction tasks, and comprising a total of 32,000 web pages from 8 different domains.

2.3 Evaluation Metrics

153

154

155

156

157

158

160

161

162

163

164

168

169

170

172

174

175

176

177

A single generation from an LLM is capable of directly extracting value from web pages and gen-

erating sequences of actions. However, the existing evaluation schemes for web page extraction tasks still follow the traditional metrics of text information extraction task, namely precision, recall, and F1 score. They limit the assessment of methods of the crawler generation task from two aspects. First, it focuses on extraction with a single webpage, rather than considering the generalizability from the perspective of a collection of webpages. Second, it does not effectively measure the transferability when adopting the action sequence to other webpages.

178

179

180

181

182

183

185

186

187

188

189

190

191

192

194

195

196

197

198

199

201

202

203

205

To address this issue, we transform the traditional IE task evaluation to executable evaluation. According to the traditional IE evaluation on a collection of webpages, we categorize the executability of action sequences into the following six situations. Specifically, for each extraction task in a websites, the result is classified based on the extraction result on precision, recall and f1-score.

(1) **Correct**, both precision, recall and f1-score equal 1, which indicates the action sequence is precisely; (2) **Precision(Prec.)**, only precision equals 1, which indicates perfect accuracy in the instances extracted following the action sequence, but miss relevant instances; (3) **Recall(Reca.)**, only recall equals 1, which means that it successfully identify all relevant instances in the webpage but incor-

¹Further details about the prompt is in Appendix B.1

²https://global.espn.com/nba/

rectly identify some irrelevant instances; (4) Unexecutable(Unex.), recall equals 0, which indicates that the action sequence fail to identify relevant instances; (5) Over-estimate(Over.), precision equals 0, which indicates that the action sequence extract the instances while ground-truth is
empty; (6) Else: the rest of the situation, including
partially extract the information, etc.

Since the above classifications are mutually exclusive, we use the ratio metric to calculate the proportion of each result in our task.

$$M_R = \frac{\# \text{ case of situation}}{\# \text{ total case}} \tag{1}$$

We are more concerned with success rate, so for the *Correct* metric, higher values indicate a better proportion of generated execution paths; whereas for the *Un-executable* metric, lower values are preferable.

3 AUTOCRAWLER

214

215

216

217

218

219

220

221

223

224

234

235

240

241

In this section, we describe our framework AU-TOCRAWLER for generating crawler to extract the specific information from semi-structure HTML. Our approach is divided into two phases: first, we adopt a progressive generation framework that utilizes the hierarchical structure of web pages; second, we employ a synthesis framework based on results from multiple web pages. The overall framework is presented in Figure 2.

3.1 Modeling

Different from the wrapper method that generate a XPath, we model the crawler generation task as a action sequence generation task. In specific, we generate an action sequence \mathcal{A}_{seq} that consists of a sequence of XPath³ expression from a set of seed webpages (i.e. a small portion of webpages in test case for generating the sequence).

 $\mathcal{A}_{seq} = [\text{XPath}_1, \text{XPath}_2, ..., \text{XPath}_n] \quad (2)$

242where n denotes the length of the action sequence.243We execute the XPath in the sequence using the244parser in order. In the sequence, all XPath expres-245sions except the last one are used for pruning the246webpage, and the last one is used for extracting the247corresponding element value from the pruned web248page.

Algorithm 1: Algorithm for progressive
understanding
Data: origin HTML code h_0 , task
instruction I , max retry times d_{max}
Result: Executable action sequence A_{seq} to
extract the value in the HTML
1 Initial history $\mathcal{A}_{seq} \leftarrow [], k = 0;$
2 while True do
// Terminate on reaching maximum steps.
3 if $k > d_{max}$ then break;
// Top-down
4 $value, xpath \leftarrow \mathbf{LLM}_g(h_k, I);$
// Get the text of the DOM node
s result \leftarrow Parser _{text} $(h_k, xpath);$
6 if <i>result</i> == <i>value</i> then break;
// Step-back
7 repeat
s $xpath \leftarrow xpath + "/";$
// Get the sub HTML of the DOM node
9 $h_{k+1} \leftarrow \mathbf{Parser}_{node}(h_k, xpath);$
10 until <i>h</i> contains value;
11 Append($\mathcal{A}_{seq}, xpath$);
12 $k \leftarrow k+1;$
13 end
14 return \mathcal{A}_{seq}

3.2 Progressive generation

Dealing with the lengthy content and hierarchical structure of web pages, generating a complete and executable web scraping path in one go encounters several key challenges: 249

250

251

252

253

254

255

256

257

258

259

261

262

263

264

265

266

267

268

269

270

271

272

The HTML content is organized in a DOM tree structure, which make it possible to prune irrelevant page components and hence, limited the length and height of the DOM tree. Specifically, we perform a traversal strategy consists of **top-down** and **stepback** operations. Top-down refers to starting from the root node of current DOM tree, progressively refine down to the specific node containing the target information. Step-back refers to the process of reassessing and adjusting selection criteria by moving up the DOM tree to choose a more reliable and broadly applicable node as a foundation for more consistent and accurate XPath targeting.

At each step, we first employ a top-down operation, guiding the LLMs to directly write out the xpath leading to the node containing the target information and judge whether the value extracted with xpath is consistent with the If execution fails, then adopt a step-back operation to retreat from

³https://en.wikipedia.org/wiki/XPath

Models	Method	EXECUTABLE EVALUATION						IE EVALUATION		
mouchs	memou	Correct(↑)	Prec	Reca	Unex.(\downarrow)	Over.	Else	Prec	Reca	F1
		Cl	osed-sou	urce LL	Ms					
	COT	36.75	8.83	6.71	43.46	0.71	3.53	89.45	50.43	47.99
GPT-3.5-Turbo	Reflextion	46.29	11.66	2.83	37.10	0.71	1.41	94.67	55.85	55.10
	AUTOCRAWLER	54.84	11.83	8.96	19.35	1.08	3.94	85.85	73.34	69.20
	СОТ	29.69	10.94	7.50	47.19	1.25	3.44	81.21	45.22	41.81
Gemini Pro	Reflextion	33.12	6.56	4.06	52.50	0.63	3.12	87.45	42.75	40.88
	AUTOCRAWLER	42.81	11.87	4.69	34.38	1.25	5.00	85.70	57.54	54.91
	СОТ	61.88	12.50	7.19	14.37	0.94	3.12	87.75	79.90	76.95
GPT4	Reflextion	<u>67.50</u>	13.75	4.37	10.94	0.94	2.50	<u>93.28</u>	<u>82.76</u>	<u>82.40</u>
	AUTOCRAWLER	71.56	14.06	5.31	4.06	0.63	4.37	92.49	89.13	88.69
Open-source LLMs										
	COT	3.44	0.31	0.63	95.31	0.00	0.63	94.23	4.55	4.24
Mistral 7B	Reflextion	2.19	0.00	0.31	97.19	0.00	0.31	95.60	2.78	2.49
	AUTOCRAWLER	2.87	0.00	0.00	96.77	0.36	0.00	98.57	3.23	2.87
	СОТ	17.98	3.75	2.25	74.53	0.00	1.50	79.75	21.98	21.36
CodeLlama	Reflextion	18.08	4.80	2.95	73.06	0.00	1.11	78.96	23.26	22.44
	AUTOCRAWLER	23.99	8.12	1.48	64.94	0.00	1.48	78.59	28.70	28.41
	СОТ	28.75	8.13	4.37	57.81	0.31	0.63	89.79	38.23	37.26
Mixtral 8×7B	Reflextion	36.25	6.88	3.12	51.25	0.00	2.50	89.35	44.57	43.60
	AUTOCRAWLER	46.88	10.62	7.19	30.31	0.63	4.37	87.32	62.71	59.75
	СОТ	36.56	10.94	5.63	42.50	0.63	3.75	86.05	48.78	47.05
Deepseek-coder	Reflextion	37.19	11.25	4.06	44.69	1.25	1.56	86.41	48.28	47.08
-	AUTOCRAWLER	38.75	11.25	5.31	39.69	0.63	4.37	84.91	52.11	49.68

Table 1: The executable evaluation and IE evaluation of LLMs with three framework in xxx task. We examine 7 LLMs, including 3 closed-source LLMs and 4 open-source LLMs.

the failed node, ensuring the web page includes the
target information, which is driven by LLMs. The
detail is shown in Algorithm 1.

3.3 Synthesis

277

281

282

284

290

291

Although we gain an executable action sequence within progressive generation process, there are still differences in the specific location of the target information and the structure between different web pages. The action sequence may collect To enhance the stability of crawler generation

We randomly select N webpages from the test case as seed webpages. Then, we generate an action sequence for each of them. Subsequently, we execute multiple different action sequences to extract information from the seed webpages respectively. We collect all action sequence and their corresponding results, and then choose one action sequence that can extract all the target information in the webpages as the final action sequence.

4 Experiment

With the goal of putting AUTOCRAWLER to practical use, we investigate the following research questions: 1) How is the overall performance of AU-TOCRAWLER in comparison to the current state-ofthe-art in crawler generation task? 2) How does our progressive understanding generation framework improve performance? What is its relationship with the size of LLM? 3) Can we entirely rely on LLMs for web automation? 4) In which scenarios do current frameworks still not perform well? 292

293

294

295

297

298

299

301

302

303

304

305

306

307

308

309

310

311

4.1 Experimental settings & evaluation metrics

We conduct our experiment on various LLMs including closed-source LLMs: **GPT-3.5-Turbo** (OpenAI, 2022), **Gemini Pro**(Team et al., 2023) and **GPT4** (OpenAI, 2023) as well as open-source LLMs: **Mistral-7B-Instruct-v0.2** (Jiang et al., 2023), **CodeLlama-34B-Instruct** (Rozière et al., 2024), **Mixtral** 8×7B-Instruct-v0.1 (Jiang et al., 2024) and Deepseek-Coder-33B-Instruct (Guo et al., 2024). Furthermore, we apply different LLM-prompt-based-web agents as our baselines, including ZS_COT (Wei et al., 2023) and Reflextion (Shinn et al., 2023) and AUTOCRAWLER to them. Due to the limited length context of LLMs, all experiments are conducted under zero-shot settings. The full prompts can be found in Appendix B.2.

Beside the execution evaluation metrics described in Section 2.3, we also employ traditional evaluation metrics to more comprehensively assess the quality of different action sequence. Specifically, We adopt precision(P.), recall(R.), and macrof1(F1), which are calculated as the mean of the corresponding metrics for each case.

4.2 Main Results

312

313

314

315

316

317

318

319

321

322

323

324

326

329

331

332

333

337

339

340

341

342

343

344

346

347

350

351

354

359

Results in Table 1 show that: 1) With AU-TOCRAWLER generating action sequence, LLMs can achieve state-of-the-art performance. Compared to the COT and Reflextion baseline, our method perform higher ratio of correct. Besides, it is worth noticing that Mixtral $8 \times 7B$ + AUTOCRAWLER can out-perform ChatGPT + Reflextion, indicating the superiority of AU-TOCRAWLER in generating executable action sequence on crawler generation task. 2) Models with small parameter sizes have significant difficulties in understanding and writing executable paths, so they can be considered challenging to apply in this task. On the contrary, larger-scale models demonstrate a more stable ability in instruction alignment, web structure comprehension, and reflection on execution result; 3) Traditional IE evaluation metrics cannot well describe the success rate of our task. Especially for the precision metric, it fails to reveal the performance gap among different methods with different models. This is because the extraction metrics only evaluate the results that have been extracted, ignoring that unexecutable or empty extractions also greatly damage the executability.

4.3 Generate with Golden Label

To better illustrate the effectiveness of our framework on generating executable action sequence, we compare the performance of the COT, Reflextion and our framework, while giving the golden label of the instruction. Through offering the same extraction targets, we can effectively detect the actual effects of different frameworks on generating action sequences.

Models	Method	EXECUTABLE EVALUATION						
	method	Correct	Prec	Reca	Unex.	Over.	Else	
Closed-source LLMs								
	COT	41.70	12.92	7.38	35.42	0.74	1.85	
GP1-3.5-	Reflextion	47.23	16.24	2.21	33.21	0.37	0.74	
Turbo	AUTOCRAWLER	56.89	19.43	5.65	13.43	0.71	3.89	
	COT	33.44	9.38	9.06	44.69	0.94	2.50	
Gemini	Reflextion	35.31	9.38	6.88	43.75	1.56	3.12	
Pro	AUTOCRAWLER	45.31	13.44	6.25	30.31	1.25	3.44	
	COT	61.88	11.56	9.06	11.56	1.25	4.69	
GPT4	Reflextion	71.25	7.19	4.69	14.37	0.94	1.56	
	AUTOCRAWLER	75.31	10.94	4.37	4.06	0.63	4.69	
	Open-source LLMs							
	COT	2.19	0.00	0.31	97.19	0.00	0.31	
Mistral 7B	Reflextion	2.19	0.00	0.00	97.50	0.31	0.00	
	AUTOCRAWLER	2.19	0.00	0.00	97.19	0.31	0.31	
	COT	21.40	6.27	2.21	66.79	0.74	2.58	
CodeLlama	Reflextion	22.21	4.93	3.94	66.95	0.49	1.48	
	AUTOCRAWLER	26.20	12.55	5.54	53.51	0.00	2.21	
Mixtral	COT	27.50	7.50	5.31	56.87	0.94	1.87	
	Reflextion	34.69	8.13	5.31	49.06	0.63	2.19	
8×7D	AUTOCRAWLER	45.62	11.56	5.94	32.50	1.25	3.12	
Deenseek	COT	35.00	18.75	5.31	36.25	0.63	4.06	
coder	Reflextion	38.75	11.87	2.81	42.19	0.63	3.75	
coder	AUTOCRAWLER	38.44	20.94	4.06	31.56	0.94	6.56	

Table 2: The executable evaluation with 7 LLMs with golden label.

Models	1	2	3	4	5	Avg.
GPT4	214	61	13	18	10	1.57
GPT-3.5-Turbo	115	65	22	30	43	2.35
Gemini Pro	94	52	33	27	105	2.99
Mixtral 8×7B	89	53	43	24	104	3.00
Mistral 7B	28	7	11	7	84	3.82
Deepseek-coder	137	70	55	29	23	2.14
CodeLlama	75	35	32	18	80	2.97

Table 3: Length of action sequence of AUTOCRAWLER based on different LLMs.

Table 2 shows experimental results, from which we can have the following observations: 1) Our proposed progressive understanding framework still effectively enhances the model's performance under this setting; 2) LLMs still suffer in accurately understanding web page contents with semistructured markup languages, which illustrate the performance gap between Table 1 and Table 2; 3) Compared to closed-source LLMs, even provided with golden labels, Open-source LLMs are unable to achieve sustained performance improvement. This phenomenon demonstrates the bottleneck for these models lies not in understanding the webpage content, but in understanding the webpage hierarchical structure itself.

362

363

364

365

366

367

368

369

370

371

372

373

374

375



Figure 3: The curve on length and height compression ratio.

4.4 Further Study with AUTOCRAWLER

379

381

385

397

398

400

401

402

403

404

405

406

The length of action sequence is dependent on the LLMs capability. In order to comprehensively explore the performance of different LLMs in understanding web page structure, we explore the impact of models on the number distribution of the steps. In specific, we collect all action sequence and calculate the average steps of AUTOCRAWLER with different LLMs under the settings describe in Section 4.1. The experimental result is reported in Table 3.

We observe that, AUTOCRAWLER with stronger LLMs generate fewer length of action sequence. AUTOCRAWLER with GPT4 generate 1.57 steps on average, while the AUTOCRAWLER with Mistral 7B generate 3.82 steps on average. This phenomenon can be interpreted as more powerful models having a better understanding of the web page hierarchical structure, thus being able to accurately output the appropriate XPaths in longer/deeper web pages, thereby reducing the number of steps.

The "U" curve of compression ratio We define the length of HTML as the number of tokens in the HTML, and its height as the height of the DOM tree represented by the HTML. we define the compression ratio of length and height as the ratio of the length/height of the original web page to that of the web page after being pruned by AUTOCRAWLER.

 $Compression_{L} = \frac{\#\text{tokens of new HTML}}{\#\text{tokens of origin HTML}}$ $Compression_{H} = \frac{\#\text{height of new HTML}}{\#\text{height of origin HTML}}$ (3)

We calculate their compression ratio of the Cor-

rect case and rank LLMs based on their performance. Figure 3 shows the result. It is interesting to note that there is a "U" curve on both the compression ratio of length and height. This phenomenon can be explained from two aspects: on one hand, when LLM is powerful, it can generate the correct XPath without the process of step-back to re-accessing the sub DOM tree; on the other hand, when the model is weak, it is unable to effectively understand the hierachical structure of web page, and thus cannot generate reliable, effective XPaths for the web page.

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

XPath fragility within AUTOCRAWLER The fragility of XPath often refers to the characteristic of XPath expressions becoming ineffective or inaccurately matching the target element when faced with new webpages. This is mainly caused by XPath specifying specific information through *predicates*, such as text, @class, etc.

We mainly focus on fragility from text because these webpages are from the same websites (i.e. @class is a good characteristic for generating stable action sequence). Table 4 show XPath expressions that rely on text. We aim to explore the reusability of generating XPath based on text features. We manually calculated the proportion of bad cases with two types of predicates, *contains* and *equal*⁴. Results in Table 5 show that the stronger LLMs capability, the lower proportion of bad case with AUTOCRAWLER. However, it is worth noticing that current SoTA LLM GPT-4 still suffers from XPath fragility problem, which indicates that relying entirely on LLMs to generate reliable XPath still has some distance to go.

4.5 Error Analysis

We perform an analysis by looking into the recorded action sequence from the AU-TOCRAWLER with GPT-4 and identify the following common failure modes. We mainly focus on the case that is categorized as unexecutable, over-estimate and else.

Non-generalizability of webpages The target information and corresponding webpage structures exhibit variations across different webpages, leading to a lack of generalizability in AUTOCRAWLER (i.e., the inability to apply the same rules across all webpages in the same website) For instance, for

⁴https://www.w3schools.com/xml/xpath_ syntax.asp

	Good case	Bad case
Question	Here's a webpage on detail information with detail information of an NBA player. Please extract the height of the player.	Here's a webpage on detail information of a university. Please extract the contact phone <i>number</i> of the university.
Case	<pre>//div[@class='gray200B-dyContent']/ b[contains(text(),'Height:')]/following- sibling::text()</pre>	//div[@class='infopage']//h5[

Table 4: Examples of XPath fragility. The green focus on the common information across different webpages, while the red focus on specific information of seed webpages.

Models	Contains	Equal(=)
GPT4	0.61%	2.90%
GPT-3.5-Turbo	9.33%	9.78%
Gemini Pro	10.62%	14.29%
Mixtral 8×7B	12.88%	8.55%
Deepseek-Coder	11.63%	7.55%
CodeLlama	18.75%	14.29%
Mistral 7B	18.18%	33.33%

Table 5: Bad case ratio in two types of predicate.

the task "Please extract the name of company that offers the job" in website job-careerbuilder, most webpages contain the company name, but there is one webpage where the company name is listed as "Not Available" on another node of DOM tree.

Miss in Multi-valued Interesting Presented with the task of generating crawler for extracting *address* in restaurant webpages or *contact phone number* in university webpages, target information is located in multiple locations in the webpage, such as the information bar, title, etc. While AU-TOCRAWLER is capable of generating action sequences to extract portions of information, crafting a comprehensive action sequence that captures all the information remains a challenge.

5 Related Work

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

5.1 Web Automation with LLMs

Many studies explore the concept of an open-world 471 in web simulation environments (Shi et al., 2017; 472 Yao et al., 2023; Deng et al., 2023; Zhou et al., 473 2023), encompassing a broad spectrum of tasks 474 found in real-life scenarios, such as online shop-475 ping, flight booking, and software development. 476 Current web automation frameworks mainly aim 477 to streamline the web environment (Sridhar et al., 478

2023; Gur et al., 2023; Zheng et al., 2024) and to devise strategies for planning and interacting with the web (Sodhi et al., 2023; Ma et al., 2023). However, these frameworks exhibit a lack of reusability, with agents heavily reliant on LLMs for even similar tasks, resulting in inefficiencies. 479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

5.2 DOM-based Web Extraction

These methods utilize the hierarchical structure of the webpage. Method of this category includes rulebased (Zheng et al., 2008), learning wrappers(i.e a DOM-specific parser that can extract content) (Gulhane et al., 2011; Kushmerick, 1997; Dalvi et al., 2011). These methods demand substantial human involvement, including the creation of wrapper annotations, the application of heuristic scoring rules (such as visual proximity), the crafting of features for neural network input, and the use of prior knowledge for verification. Contemporary strategies employ distant supervision to autonomously create training samples by matching data from existing knowledge bases (KBs) with web sources (Lockard et al., 2018, 2019). While this significantly lowers the effort required for annotation, it unavoidably leads to false negatives because of the incompleteness of knowledge bases (KBs) (Xie et al., 2021).

6 Conclusion

In this paper, we introduce the crawler generation task and the paradigm that combines LLMs and crawlers together to improve the reusability of the current web automation framework. We then propose AUTOCRAWLER, a two-phase progressive understanding framework to generate a more stable and executable action sequence. Our comprehensive experiments demonstrate that AU-TOCRAWLER can outperform the state-of-the-art baseline in the crawler generation task.

- 566 567
- 568
- 569 570
- 571 572 573 574
- 575
- 576 577 578 579 580

- 585 586 587 588
- 589 590
- 591 592 593
- 594 595 596
- 597 598
- 599 600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

Limitation 515

525

527

529

532

533

534

537

539

545

547

553

554

555

556

557

516 We introduce a paradigm that combine large language models (LLMs) with crawler for web crawler 517 generation task, and propose AUTOCRAWLER to 518 generate a executable action sequence with pro-519 gressively understanding the HTML documents. 521 Though experimental results show the effectiveness of our framework, there are still some limits of our works. 523

> First, our framework is restricted with the paradigm in vertical webpages information extraction task. LLMs with crawler provide high efficiency on open-world web IE task, but it can hardly transfer to existing web environment such as Mind2Web (Deng et al., 2023), WebArena (Zhou et al., 2023).

Second, the testbed dataset (only SWDE (Hao et al., 2011) dataset) is limited because there is currently a lack of data sets related to web page research. Building a comprehensive and widespread dataset for the web crawler generation task is our future work.

Ethic statement

We hereby declare that all authors of this article are aware of and adhere to the provided ACL Code of Ethics and honor the code of conduct.

Use of Human Annotations Human annotations 541 are only utilized in the early stages of methodological research to assess the feasibility of the proposed solution. All annotators have provided consent for 544 the use of their data for research purposes. We guarantee the security of all annotators throughout 546 the annotation process, and they are justly remunerated according to local standards. Human annotations are not employed during the evaluation of our 549 method.

> **Risks** The datasets used in the paper have been obtained from public sources and anonymized to protect against any offensive information. Though we have taken measures to do so, we cannot guarantee that the datasets do not contain any socially harmful or toxic language.

References

Mirko Bronzi, Valter Crescenzi, Paolo Merialdo, and Paolo Papotti. 2013. Extraction and integration of partially overlapping web sources. Proc. VLDB Endow., 6(10):805-816.

- Nilesh Dalvi, Ravi Kumar, and Mohamed Soliman. 2011. Automatic wrappers for large scale web extraction. arXiv preprint arXiv:1103.2406.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web.
- Pankaj Gulhane, Amit Madaan, Rupesh Mehta, Jeyashankher Ramamirtham, Rajeev Rastogi, Sandeep Satpal, Srinivasan H Sengamedu, Ashwin Tengli, and Charu Tiwari. 2011. Web-scale information extraction with vertex. In 2011 IEEE 27th International Conference on Data Engineering, pages 1209-1220.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. Deepseek-coder: When the large language model meets programming – the rise of code intelligence.
- Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. 2023. A real-world webagent with planning, long context understanding, and program synthesis.
- Qiang Hao, Rui Cai, Yanwei Pang, and Lei Zhang. 2011. From one tree to a forest: a unified solution for structured web data extraction. In *Proceedings of the 34th* International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '11, page 775–784, New York, NY, USA. Association for Computing Machinery.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts.
- Nicholas Kushmerick. 1997. Wrapper induction for information extraction. University of Washington.
- Colin Lockard, Xin Luna Dong, Arash Einolghozati, and Prashant Shiralkar. 2018. Ceres: Distantly supervised relation extraction from the semi-structured web. arXiv preprint arXiv:1804.04635.

Colin Lockard, Prashant Shiralkar, and Xin Luna Dong. 2019. Openceres: When open information extraction meets the semi-structured web. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 3047–3056.

618

619

635

642

643

645

646

647

651 652

655

667 668

672

- Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, and Dong Yu. 2023. Laser: Llm agent with state-space exploration for web navigation.
- OpenAI. 2022. https://openai.com/ chatgpt/. Accessed: 2024-02-07.
 - OpenAI. 2023. https://openai.com/gpt-4/. Accessed: 2024-02-07.
 - Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code llama: Open foundation models for code.
 - Tianlin Shi, Andrej Karpathy, Linxi (Jim) Fan, Josefa Z. Hernández, and Percy Liang. 2017. World of bits: An open-domain platform for web-based agents. In International Conference on Machine Learning.
 - Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao.
 2023. Reflexion: Language agents with verbal reinforcement learning.
 - Paloma Sodhi, S. R. K. Branavan, and Ryan McDonald. 2023. Heap: Hierarchical policies for web actions using llms.
 - Abishek Sridhar, Robert Lo, Frank F. Xu, Hao Zhu, and Shuyan Zhou. 2023. Hierarchical prompting assists large language model on web navigation.
 - Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, and Anja Hauth. 2023. Gemini: A family of highly capable multimodal models.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models.
 - Chenhao Xie, Jiaqing Liang, Jingping Liu, Chengsong Huang, Wenhao Huang, and Yanghua Xiao. 2021.
 Revisiting the negative data of distantly supervised relation extraction. *arXiv preprint arXiv:2105.10158*.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2023. Webshop: Towards scalable real-world web interaction with grounded language agents.

Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. 2024. Synapse: Trajectory-as-exemplar prompting with memory for computer control. 673

674

675

676

677

678

679

680

681

682

683

684

- Xiaolin Zheng, Tao Zhou, Zukun Yu, and Deren Chen. 2008. Url rule based focused crawler. In 2008 IEEE International Conference on e-Business Engineering, pages 147–154. IEEE.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. 2023. Webarena: A realistic web environment for building autonomous agents.

- 68
- 687

....

- 68
- 690
- 69
- 69
- 69
- 69
- 696
- 697

Listing 1: Prompts for ZS_COT, Reflextion and AutoCrawler

Table 6 shows the detail statistic about semi-

structure web information extraction dataset

Table 7 shows the task prompt we design for each

We provide a comprehensive list of all the prompts

that have been utilized in this study, offering a

clear reference for understanding our experimental

A Dataset Statistic

B Prompt List

B.1 Task Prompt

attribute in origin dataset.

B.2 Module Prompt

SWDE.

approach.

```
Prompt-crawler (for COT, Reflextion and
           AutoCrawler):
700
           Please read the following HTML code, and
701
            then return an Xpath that can recognize
            the element in the HTML matching the
702
703
           instruction below.
           Instruction: {0}
704
705
           Here're some hints:
           1. Do not output the xpath with exact
706
           value or element appears in the HTML.
           2. Do not output the xpath that indicate
           multi node with different value. It
710
           would be appreciate to use more @class
           to identify different node that may
711
712
           share the same xpath expression.
713
           3. If the HTML code doesn't contain the
714
           suitable information match the
715
           instruction, keep the xpath attrs blank.
716
           4. Avoid using some string function such
           as 'substring()' and 'normalize-space()
717
718
           ' to normalize the text in the node.
           Please output in the following Json
720
           format:
721
722
           {
723
               "thought": "", # a brief thought of
724
               how to confirm the value and
725
               generate the xpath
               "value": "", # the value extracted
726
727
               from the HTML that match the
728
               instruction
               "xpath": "", # the xpath to extract
729
               the value
731
           }
732
           Here's the HTML code:
           ...
733
734
           \{1\}
735
736
737
           Prompt-reflextion (for Reflextion and
738
           AutoCrawler):
739
           Here's the HTML extraction task:
740
           Task description: Please read the
741
           following HTML code, and then return an
```

Xpath that can recognize the element in the HTML matching the instruction below. Instruction: {0} We will offer some history about the thought and the extraction result. Please reflect on the history trajectory and adjust the xpath rule for better and more exact extraction. Here's some hints: 1. Judge whether the results in the history is consistent with the expected value. Please pay attention for the following case: 1) Whether the extraction result contains some elements that is irrelevent 2) Whether the crawler return a empty result 3) The raw values containing redundant separators is considered as consistent because we will postprocess it. 2. Re-thinking the expected value and how to find it depend on xpath code 3. Generate a new or keep the origin xpath depend on the judgement and thinking following the hints: 1. Do not output the xpath with exact value or element appears in the HTML. 2. Do not output the xpath that indicate multi node with different value. It would be appreciate to use more @class and [num] to identify different node that may share the same xpath expression. 3. If the HTML code doesn't contain the suitable information match the instruction, keep the xpath attrs blank. Please output in the following json format: "thought": "", # thought of why the xpaths in history are not work and how to adjust the xpath "consistent": "", # whether the extracted result is consistent with the expected value, return yes/no directly "value": "", # the value extracted from the HTML that match the task description "xpath": "", # a new xpath that is different from the xpath in the following history if not consistent } And here's the history about the thought , xpath and result extracted by crawler. $\{1\}$ Here's the HTML code: • • • {2} **Prompt**-synthesis (for COT, Reflextion and AutoCrawler):

742

743

744

745

746

747

748 749

750

751

752

753

754

755

756 757

758 759

760

761

762

763

764

765

766

767 768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

```
You're a perfect discriminator which is
812
813
           good at HTML understanding as well.
814
           Following the instruction, there are
815
           some action sequence written from
816
           several HTML and the corresponding
817
           result extracted from several HTML.
818
           Please choose one that can be best
           potentially adapted to the same
819
820
           extraction task on other webpage in the
821
           same websites. Here are the instruction
822
           of the task:
823
           Instructions: {0}
824
           The action sequences and the
825
           corresponding extracted results with
826
           different sequence on different webpage
827
           are as follow:
828
           {1}
829
830
           Please output the best action sequence
           in the following Json format:
831
832
                "thought": "" # brief thinking about
833
834
                which to choose
                "number": "" # the best action
835
836
               sequence choosen from the candidates
837
                , starts from 0. If there is none,
838
               output 0.
839
           }
840
           Prompt-judgement (for AutoCrawler):
841
842
           Your main task is to judge whether the
843
           extracted value is consistent with the
844
           expected value, which is recognized
           beforehand. Please pay attention for the
845
846
            following case:
847
               1) If the extracted result contains
848
               some elements that is not in
849
               expected value, or contains empty
               value, it is not consistent.
851
               2) The raw values containing
852
               redundant separators is considered
853
               as consistent because we can
854
               postprocess it.
855
856
           The extracted value is: {0}
857
           The expected value is: {1}
858
859
           Please output your judgement in the
           following Json format:
861
           { {
                "thought": "", # a brief thinking
862
               about whether the extracted value is
863
864
                consistent with the expected value
                "judgement": "" # return yes/no
865
866
               directly
867
           } }
868
869
           Prompt-stepback (for AutoCrawler):
870
           Your main task is to judge whether the
871
           following HTML code contains all the
872
           expected value, which is recognized
873
           beforehand.
874
           Instruction: {0}
875
           And here's the value: {1}
876
           The HTML code is as follow:
877
           ...
878
           {2}
879
           ...
880
```

Please output your judgement in the	881
following Json format:	882
{	883
"thought": "", # a brief thinking	884
about whether the HTML code contains	885
expected value	886
"judgement": "" # whether the HTML	887
code contains all extracted value.	888
Return yes/no directly.	889
}	890

Domain	Attribute	Website	Num	Domain	Attribute	Website	Num
Auto	model price engine fuel_economy	aol autobytel automotive autoweb carquotes cars kbb motortrend msn yahoo	2000 2000 1999 2000 2000 657 2000 1267 2000 2000	Movie	ttitle director genre mpaa_rating	allmovie amctv boxofficemojo hollywood iheartmovies imdb metacritic msn rottentomatoes yahoo	2000 2000 2000 2000 2000 2000 2000 200
Book	title author isbn_13 publisher pub_date	abebooks amazon barnesandnoble bookdepository booksamillion bookorders buy christianbook deepdiscount waterstone	2000 2000 2000 2000 2000 2000 2000 200	NBAPlayer	name team height weight	espn fanhouse foxsports msnca nba si slam usatoday wiki yahoo	434 446 425 434 434 515 423 436 420 438
Camera	model price manufacturer	amazon beachaudio buy compsource ecost jr newegg onsale pcnation thenerd	1767 247 500 430 923 367 220 261 234 309	Restaurant	name address phone cuisine	fodors frommers zagat gayot opentable pickaretaurant restaurantica tripadvisor urbanspoon usdiners	2000 2000 2000 2000 2000 2000 2000 200
Job	title company location date_posted	careerbuilder dice hotjobs job jobcircle jobtarget monster nettemps rightitjobs techcentric	2000 2000 2000 2000 2000 2000 2000 200	University	name phone website type	collegeboard collegenavigato collegeprowler collegetoolkit ecampustours embark matchcollege princetonreview studentaid usnews	2000 2000 2000 2000 1063 2000 2000 v 615 2000 1027

Table 6: Detail statistic of SWDE dataset.

Domain	Task prompt	Prompt
Auto	Here's a webpage with detail in- formation of an auto.	Please extract the model of the auto. Please extract the price of the auto. Please extract the engine of the auto. Please extract the fuel efficiency of the auto.
Book	Here's a webpage with detail in- formation of a book.	Please extract the title of the book. Please extract the author of the book. Please extract the isbn number of the book. Please extract the publisher of the book. Please extract the publication date of the book.
Camera	Here's a webpage with detail in- formation of camera.	Please extract the product name of the camera. Please extract the sale price of the camera. Please extract the manufactor of the camera.
Job	Here's a webpage with detail in- formation of a job.	Please extract the title of the job. Please extract the name of company that offers the job. Please extract the working location of the job. Please extract the date that post the job.
Movie	Here's a webpage with detail in- formation of a movie.	Please extract the title of the movie. Please extract the director of the movie. Please extract the genre of the movie. Please extract the MPAA rating of the movie.
NBAPlayer	Here's a webpage with detail in- formation of an NBA player.	Please extract the name of the player. Please extract the team of the player he play now. Please extract the height of the player. Please extract the weight of the player.
Restaurant	Here's a webpage with detail in- formation of a restaurant.	Please extract the restaurant's name. Please extract the retaurant's address. Please extract the restaurant's phone number. Please extract the cuisine that the restaurant offers.
University	Here's a webpage on detail in- formation of a university.	Please extract the name of the university. Please extract the contact phone number of the university. Please extract the website url of the university. Please extract the type of the university.

Table 7: Prompts for crawler generation task in this paper.