

Introduction to WebQL

A Modern Web Data Extraction Language

Version 1.0

Ylli Prifti
ylli@prifti.us

September 12, 2025

Abstract

WebQL (Web Query Language) is a declarative language for web data extraction that combines the simplicity of JSON5 syntax with powerful extraction capabilities. Building upon the foundations of OXPath and modern web technologies, WebQL provides a unified approach to extracting structured data from web pages, handling dynamic content, multi-modal data, and integrating AI-powered selectors. This manual provides a comprehensive introduction to WebQL, its syntax, semantics, and practical applications in modern web data extraction scenarios.

Contents

1	Introduction	4
1.1	Motivation	4
1.2	Design Principles	4
1.3	Comparison with Existing Approaches	4
2	Getting Started	5
2.1	Installation	5
2.2	Your First Query	5
2.3	Understanding the Structure	5
3	Language Syntax	5
3.1	Basic Structure	5
3.2	Step Types	6
3.2.1	Extract Step	6
3.2.2	Action Steps	6
3.2.3	Control Flow Steps	6
3.3	Selectors	7
4	Advanced Features	7
4.1	Pagination Handling	7
4.2	Dynamic Content	7
4.3	Multi-modal Extraction	8
4.4	AI-Powered Extraction	8
5	Plugin System	8
5.1	Built-in Plugins	8
5.2	Using Plugins in Queries	9
5.3	Creating Custom Plugins	9
6	Query Patterns	9
6.1	E-commerce Product Extraction	9
6.2	News Article Extraction	10
6.3	Form Interaction and Search	10
7	Error Handling and Debugging	11
7.1	Error Handling Strategies	11
7.2	Debugging Queries	11
7.3	Common Issues and Solutions	12
8	Performance Optimization	12
8.1	Parallel Extraction	12
8.2	Caching	12
8.3	Selective Extraction	12
9	Integration	13
9.1	Python Integration	13
9.2	REST API	13
9.3	Scheduling and Automation	13
10	Best Practices	13
10.1	Query Design	13

10.2 Reliability	14
10.3 Ethics and Compliance	14
11 Case Studies	14
11.1 Academic Research: Citation Network Analysis	14
11.2 Real Estate Market Analysis	15
12 Conclusion	15
12.1 Future Directions	15
12.2 Community and Resources	16
A Appendix A: Complete Language Reference	16
A.1 Query Structure	16
A.2 Step Types Reference	16
B Appendix B: Selector Reference	16
C Appendix C: Plugin API	16

1 Introduction

The web has become the largest repository of information in human history, yet extracting structured data from web pages remains a challenging task. Traditional approaches using imperative programming with libraries like BeautifulSoup or Puppeteer require significant coding effort and domain expertise. WebQL addresses this challenge by providing a declarative language specifically designed for web data extraction.

1.1 Motivation

Modern web applications present unique challenges for data extraction:

- **Dynamic Content:** Single-page applications (SPAs) load content dynamically
- **Multi-modal Data:** Web pages contain text, images, videos, and embedded documents
- **Complex Interactions:** Data often requires user interactions (clicks, scrolls, form submissions)
- **Anti-bot Measures:** Websites implement CAPTCHAs and rate limiting
- **Maintenance Burden:** Selectors break when websites update their structure

WebQL addresses these challenges through:

- Declarative syntax that focuses on *what* to extract, not *how*
- Built-in support for browser automation and interactions
- AI-powered selectors that adapt to structural changes
- Plugin architecture for extending functionality
- Native handling of multi-modal content

1.2 Design Principles

WebQL is designed around four core principles:

1. **Simplicity:** Use familiar JSON5 syntax that's readable and writable by humans
2. **Composability:** Build complex extractions from simple, reusable components
3. **Extensibility:** Support custom processors through a plugin architecture
4. **Robustness:** Handle failures gracefully with built-in retry and error handling

1.3 Comparison with Existing Approaches

Feature	WebQL	OXPath	Scrapy	Puppeteer
Declarative	✓	✓	Partial	✗
Browser Automation	✓	✓	✗	✓
AI Integration	✓	✗	✗	✗
Multi-modal	✓	✗	Partial	Partial
Plugin System	✓	✗	✓	✗
JSON Output	Native	Extension	✓	Manual

Table 1: Comparison of WebQL with existing web extraction tools

2 Getting Started

2.1 Installation

WebQL is implemented through the DR Web Engine, which can be installed via pip:

```
pip install dr-web-engine
playwright install chromium
```

2.2 Your First Query

Let's start with a simple example that extracts quotes from a website:

```
1  {
2      "@url": "https://quotes.toscrape.com",
3      "@steps": [
4          {
5              "@xpath": "//div[@class='quote'][1]",
6              "@fields": {
7                  "text": "./span[@class='text']/text()", // Relative XPath
8                  "author": "./small[@class='author']/text()" // Relative XPath
9              }
10         }
11     ]
12 }
```

To execute this query:

```
drweb -q quotes.json5 -o results.json
```

2.3 Understanding the Structure

Every WebQL query consists of:

- **@url**: The target URL for extraction
- **@steps**: An array of extraction steps, each containing:
 - **@xpath**: XPath selector to find elements
 - **@fields**: Field mappings using relative XPath
 - Optional format directives like **@format**, **@output**

3 Language Syntax

3.1 Basic Structure

A WebQL query follows this structure:

```
1  {
2      "@url": "https://example.com", // Required: Target URL
3      "@steps": [ // Required: Extraction steps
4          {
5              "@xpath": "//div[@class='item']", // XPath selector
6              "@fields": { // Field mappings
7                  "title": "./h2/text()", // Relative XPath
8                  "description": "./p/text()"
9              }
10         }
11     ]
12 }
```

3.2 Core Directives

WebQL uses directives prefixed with @ to define extraction behavior:

3.2.1 Basic Extraction

The fundamental extraction step uses @xpath and @fields:

```

1 {
2   "@xpath": "//div[@class='quote']",
3   "@fields": {
4     "text": "./span[@class='text']/text()", 
5     "author": "./small[@class='author']/text()"
6   }
7 }
```

3.2.2 Output Format Directives

Control output format and structure:

```

1 {
2   "@format": "jsonl",                                // Output format
3   "@output": "results.jsonl",                         // Output filename
4   "@compression": "minimal",                          // Compression level
5   "@include-metadata": true,                          // Include extraction metadata
6   "@streaming": true                                // Stream results
7 }
```

3.2.3 LLM Training Format

Format data for AI model training:

```

1 {
2   "@format": "openai-chat",
3   "@system-prompt": "You are a helpful assistant.",
4   "@user-template": "Quote: {text} by {author}",
5   "@field-mapping": {
6     "quote_text": "text",
7     "quote_author": "author"
8   }
9 }
```

3.2.4 Navigation and Following Links

Follow links to extract data from multiple pages:

```

1 {
2   "@follow": {
3     "@xpath": "./a[contains(@href, '/details/')]",
4     "@max-depth": 2,
5     "@steps": [
6       {
7         "@xpath": "//div[@class='detail-content']",
8         "@fields": {
9           "detailed_info": "./p/text()", 
10          "specs": "./ul[@class='specs']/li/text()"
11        }
12      }
13    ]
14  }
15 }
```

3.2.5 Nested Step Definitions

Use `@step` to define nested extraction logic:

```

1  {
2      "@format": "jsonl",
3      "@step": {
4          "@xpath": "//article[@class='news-item']",
5          "@fields": {
6              "headline": "./h2/text()",  

7              "summary": "./p[@class='summary']/text()",  

8              "date": "./time/@datetime"
9          }
10     }
11 }
```

3.3 XPath Selectors

WebQL primarily uses XPath for element selection:

- **Element Selection:** `"//div[@class='content']"`
- **Text Extraction:** `".//span/text()"`
- **Attribute Values:** `".//a/@href"`
- **Relative Paths:** Use `".//"` for relative to current context
- **Multiple Matches:** XPath naturally handles multiple matching elements

4 Advanced Features

4.1 Output Formatting for LLM Training

WebQL supports various output formats optimized for machine learning:

```

1  {
2      "@url": "https://quotes.toscrape.com",
3      "@steps": [
4          {
5              "@format": "openai-chat",
6              "@output": "quotes_training.jsonl",
7              "@system-prompt": "You are a wisdom analyzer.",
8              "@user-template": "Quote: {text} by {author}",
9              "@compression": "minimal",
10             "@include-metadata": true,
11             "@step": {
12                 "@xpath": "//div[@class='quote'][1]",
13                 "@fields": {
14                     "text": "./span[@class='text']/text()",  

15                     "author": "./small[@class='author']/text()"
16                 }
17             }
18         }
19     ]
20 }
```

4.2 Multi-Page Navigation

Follow links to extract data from multiple pages:

```

1  {
2      "@curl": "https://example.com/articles",
3      "@steps": [
4          {
5              "@xpath": "//article[@class='summary']",
6              "@fields": {
7                  "title": "./h2/text()",
8                  "summary": "./p/text()"
9              },
10             "@follow": {
11                 "@xpath": "./a[@class='read-more']",
12                 "@max-depth": 2,
13                 "@steps": [
14                     {
15                         "@xpath": "//article[@class='full-content']",
16                         "@fields": {
17                             "full_content": "./div[@class='content']//text()",  

18                             "tags": "./span[@class='tag']/text()"
19                         }
20                     }
21                 ]
22             }
23         }
24     ]
25 }
```

4.3 Multi-modal Extraction

Extract different types of content:

```

1  {
2      "@extract": {
3          text_content: ".article-body",
4          images: [
5              "@base": "img",
6              url: "@attr:src",
7              alt: "@attr:alt",
8              ocr_text: "@ocr" // Extract text from image
9          ],
10         pdf_link: {
11             selector: "a[href$=.pdf]",
12             url: "@attr:href",
13             content: "@pdf_text" // Extract text from PDF
14         }
15     }
16 }
```

4.4 AI-Powered Extraction

Use natural language to describe what to extract:

```

1  {
2      "@ai_extract": {
3          prompt: "Extract product specifications as key-value pairs",
4          selector: ".specifications-section"
5      }
6  }
7 }
```

```

8 {
9   "@ai_select": {
10     description: "Find the 'Add to Cart' button",
11     action: "click"
12   }
13 }
```

5 Plugin System

5.1 Built-in Plugins

WebQL includes several built-in plugins:

- **smart-retry**: Intelligent retry with exponential backoff
- **proxy-rotation**: Automatic proxy rotation
- **ai-selector**: Natural language selectors
- **api-extractor**: API response extraction
- **jsonld-extractor**: JSON-LD structured data extraction

5.2 Using Plugins in Queries

```

1 {
2   url: "https://example.com",
3   plugins: {
4     "smart-retry": {
5       max_attempts: 3,
6       backoff: "exponential"
7     }
8   },
9   steps: [
10   {
11     "@with-retry": {
12       "@extract": {
13         data: ".fragile-selector"
14       }
15     }
16   }
17 ]
18 }
```

5.3 Creating Custom Plugins

Extend WebQL with custom functionality:

```

from engine.web_engine.plugin_interface import DrWebPlugin

class CustomPlugin(DrWebPlugin):
    @property
    def metadata(self):
        return PluginMetadata(
            name="custom-extractor",
            version="1.0.0",
            description="Custom extraction logic")
```

```

        )

def get_processors(self):
    return [CustomProcessor()]

```

6 Query Patterns

6.1 E-commerce Product Extraction

```

1 {
2     url: "https://shop.example.com/products",
3     steps: [
4         {
5             "@loop": {
6                 while: "exists(.product-card)",
7                 steps: [
8                     {
9                         "@extract": {
10                            products: [
11                                "@base": ".product-card",
12                                name: ".product-title",
13                                price: ".price-now",
14                                original_price: ".price-was",
15                                rating: ".rating-value",
16                                reviews: ".review-count",
17                                image: "img@attr:src",
18                                link: "a@attr:href"
19                            ]
20                        }
21                    },
22                    {
23                        "@click": ".load-more-products",
24                        "@wait": 2000
25                    }
26                ]
27            }
28        }
29    ]
30}

```

6.2 News Article Extraction

```

1 {
2     url: "https://news.example.com",
3     steps: [
4         {
5             "@extract": {
6                 articles: [
7                     "@base": "article",
8                     headline: "h1",
9                     author: ".author-name",
10                    date: "time@attr:datetime",
11                    content: ".article-body",
12                    tags: [".tag"],
13                    related: [
14                        "@base": ".related-article",
15                        title: ".title",
16                        url: "a@attr:href"
17                    ]
18                ]
19            }
20        }
21    ]
22}

```

```

18     }] ]
19   }
20 ]
21 }
22 }
```

6.3 Form Interaction and Search

```

1 {
2   url: "https://example.com/search",
3   steps: [
4     {
5       "@type": {
6         selector: "input[name='query']",
7         text: "machine learning"
8       }
9     },
10    {
11      "@select": {
12        selector: "select[name='category']",
13        value: "research"
14      }
15    },
16    {
17      "@click": "button[type='submit']"
18    },
19    {
20      "@wait": "selector:.search-results"
21    },
22    {
23      "@extract": {
24        results: [
25          "@base": ".result-item",
26          title: ".result-title",
27          url: "a@attr:href",
28          snippet: ".result-snippet"
29        ]
30      }
31    }
32  ]
33 }
```

7 Error Handling and Debugging

7.1 Error Handling Strategies

WebQL provides multiple error handling mechanisms:

```

1 {
2   url: "https://example.com",
3   on_error: "continue", // continue, stop, retry
4   steps: [
5     {
6       "@try": {
7         steps: [
8           {"@extract": {data: ".may-not-exist"}}
9         ],
10        catch: [
11          {"@log": "Extraction failed, using fallback"},
12          {"@extract": {data: ".fallback-selector"}}
13        ]
14      }
15    }
16  ]
17 }
```

```

13     ]
14   }
15 }
16 ]
17 }
```

7.2 Debugging Queries

Enable debug mode for detailed execution logs:

```
drweb -q query.json5 -o output.json --log-level debug
```

7.3 Common Issues and Solutions

Issue	Solution
Selector returns empty	Use browser DevTools to verify selector; try AI selector as fallback
Dynamic content not loaded	Add <code>@wait</code> step or use <code>wait_until: "networkidle"</code>
Rate limiting	Enable proxy-rotation plugin or add delays between requests
CAPTCHA blocking	Use residential proxies or implement CAPTCHA solving service

Table 2: Common issues and their solutions

8 Performance Optimization

8.1 Parallel Extraction

Execute multiple extractions concurrently:

```

1 {
2   urls: [
3     "https://example.com/page1",
4     "https://example.com/page2",
5     "https://example.com/page3"
6   ],
7   parallel: 3, // Process 3 URLs concurrently
8   steps: [
9     {"@extract": { /* ... */}}
10  ]
11 }
```

8.2 Caching

Enable caching to avoid redundant requests:

```

1 {
2   config: {
3     cache: {
4       enabled: true,
5       ttl: 3600, // Cache for 1 hour
6       storage: "redis://localhost:6379"
7     }
8   }
9 }
```

8.3 Selective Extraction

Optimize by extracting only required fields:

```

1  {
2      "@extract": {
3          // Only extract if condition is met
4          "@if": "exists(.premium-content)",
5          premium_data: ".premium-content"
6      }
7 }
```

9 Integration

9.1 Python Integration

```

from dr_web_engine import WebQL

# Initialize engine
engine = WebQL()

# Define query
query = {
    "url": "https://example.com",
    "steps": [
        {"@extract": {"title": "h1"}}
    ]
}

# Execute
results = engine.execute(query)
print(results)
```

9.2 REST API

WebQL can be exposed as a REST API:

```

# Start API server
drweb-server --port 8080

# Submit query
curl -X POST http://localhost:8080/extract \
-H "Content-Type: application/json" \
-d @query.json5
```

9.3 Scheduling and Automation

Integrate with schedulers for automated extraction:

```

# cron-job.yaml
schedule: "0 */6 * * *" # Every 6 hours
query:
    url: "https://example.com/prices"
    steps:
```

```

    - "@extract":
      prices: [".price"]
output:
  type: "webhook"
  url: "https://api.example.com/prices"

```

10 Best Practices

10.1 Query Design

1. **Start Simple:** Begin with basic selectors, add complexity gradually
2. **Use Specific Selectors:** Prefer IDs and unique classes over generic tags
3. **Test Incrementally:** Verify each step before adding the next
4. **Handle Edge Cases:** Account for missing elements and pagination limits
5. **Document Queries:** Add comments to explain complex logic

10.2 Reliability

1. **Implement Retries:** Use smart-retry plugin for transient failures
2. **Fallback Selectors:** Provide alternative selectors for critical data
3. **Monitor Changes:** Set up alerts for extraction failures
4. **Version Control:** Track query changes in Git
5. **Test Suite:** Maintain test queries for regression testing

10.3 Ethics and Compliance

1. **Respect robots.txt:** Check and comply with site policies
2. **Rate Limiting:** Implement delays to avoid overwhelming servers
3. **User-Agent:** Identify your bot appropriately
4. **Terms of Service:** Review and comply with website ToS
5. **Data Privacy:** Handle extracted data according to regulations

11 Case Studies

11.1 Academic Research: Citation Network Analysis

```

1 {
2   url: "https://scholar.google.com/scholar?q=machine+learning",
3   steps: [
4     {
5       "@follow": {
6         selector: ".gs_r",
7         max_items: 100,
8         extract_each: {
9           title: "h3 a",

```

```

10     authors: ".gs_a",
11     year: ".gs_a:regex(\\d{4})",
12     citations: "a:contains('Cited by')",
13     pdf_link: "a[href$=.pdf]@attr:href"
14   },
15   follow_next: "button[aria-label='Next']"
16 }
17 ]
18 }
19 }
```

11.2 Real Estate Market Analysis

```

1 {
2   url: "https://realestate.example.com/listings",
3   config: {
4     viewport: {width: 1920, height: 1080}
5   },
6   steps: [
7     {
8       "@interact": {
9         filters: {
10           "@type": {
11             selector: "#price-min",
12             text: "500000"
13           },
14           "@type": {
15             selector: "#price-max",
16             text: "1000000"
17           },
18           "@select": {
19             selector: "#property-type",
20             value: "house"
21           },
22           "@click": "button.apply-filters"
23         }
24       }
25     },
26     {
27       "@wait": 3000
28     },
29     {
30       "@extract": {
31         listings: [
32           "@base": ".listing-card",
33           address: ".address",
34           price: ".price",
35           bedrooms: ".beds",
36           bathrooms: ".baths",
37           sqft: ".square-feet",
38           image: "img@attr:src",
39           description: ".description",
40           agent: ".agent-name",
41           listed_date: ".list-date"
42         ]
43       }
44     ]
45   }
46 }
```

12 Conclusion

WebQL represents a significant advancement in web data extraction technology, combining declarative simplicity with powerful extraction capabilities. By abstracting away the complexities of web scraping while providing fine-grained control when needed, WebQL enables both technical and non-technical users to extract structured data from the modern web efficiently.

12.1 Future Directions

The WebQL language continues to evolve with planned enhancements including:

- **Visual Query Builder:** GUI for creating queries without writing code
- **ML-based Selector Generation:** Automatic selector learning from examples
- **Distributed Execution:** Built-in support for distributed scraping
- **Schema Validation:** Type system for extracted data validation
- **Query Optimization:** Automatic query optimization for performance

12.2 Community and Resources

- **Website:** <https://webql.dev>
- **GitHub:** <https://github.com/starlitlog/dr-web-engine>
- **Documentation:** <https://docs.webql.dev>
- **Community Forum:** <https://discuss.webql.dev>

A Appendix A: Complete Language Reference

A.1 Query Structure

```

1 {
2   // Required fields
3   url: String | Array<String>,
4   steps: Array<Step>,
5
6   // Optional fields
7   config: Config,
8   plugins: PluginConfig,
9   on_error: "stop" | "continue" | "retry",
10  parallel: Number,
11  cache: CacheConfig
12 }
```

A.2 Step Types Reference

[Complete reference of all step types with parameters...]

B Appendix B: Selector Reference

[Complete list of selector types and syntax...]

C Appendix C: Plugin API

[Plugin development guide and API reference...]