

# Atra

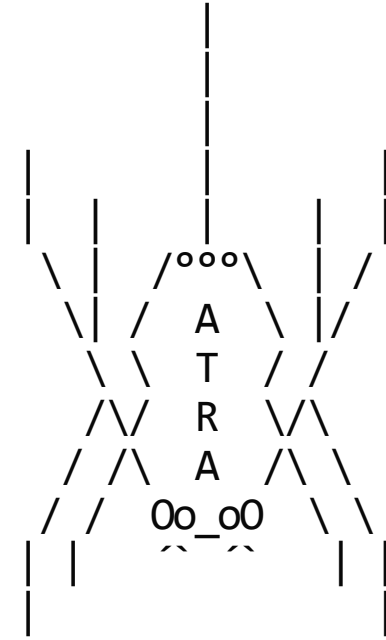
A powerful, lightweight approach to crawling



Universität Bamberg  
Lehrstuhl für Medieninformatik

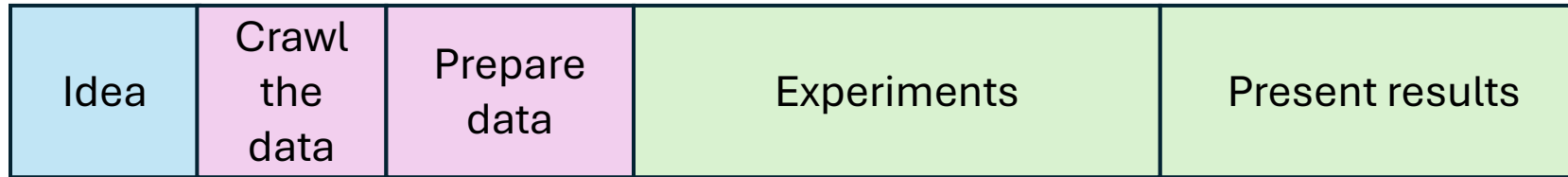
# Atra

- Motivation
- Key Features
- Roadmap
- Demo



# Motivation

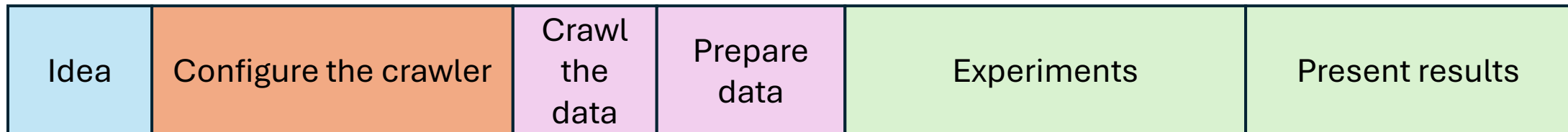
## Perfect World



## The usual for experts

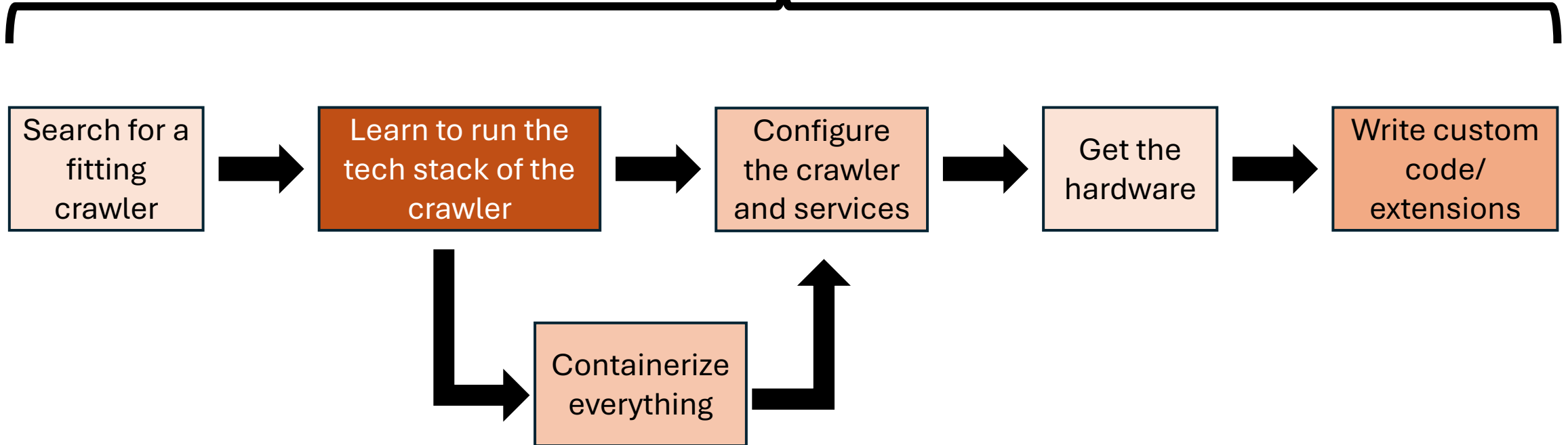


## New to crawling

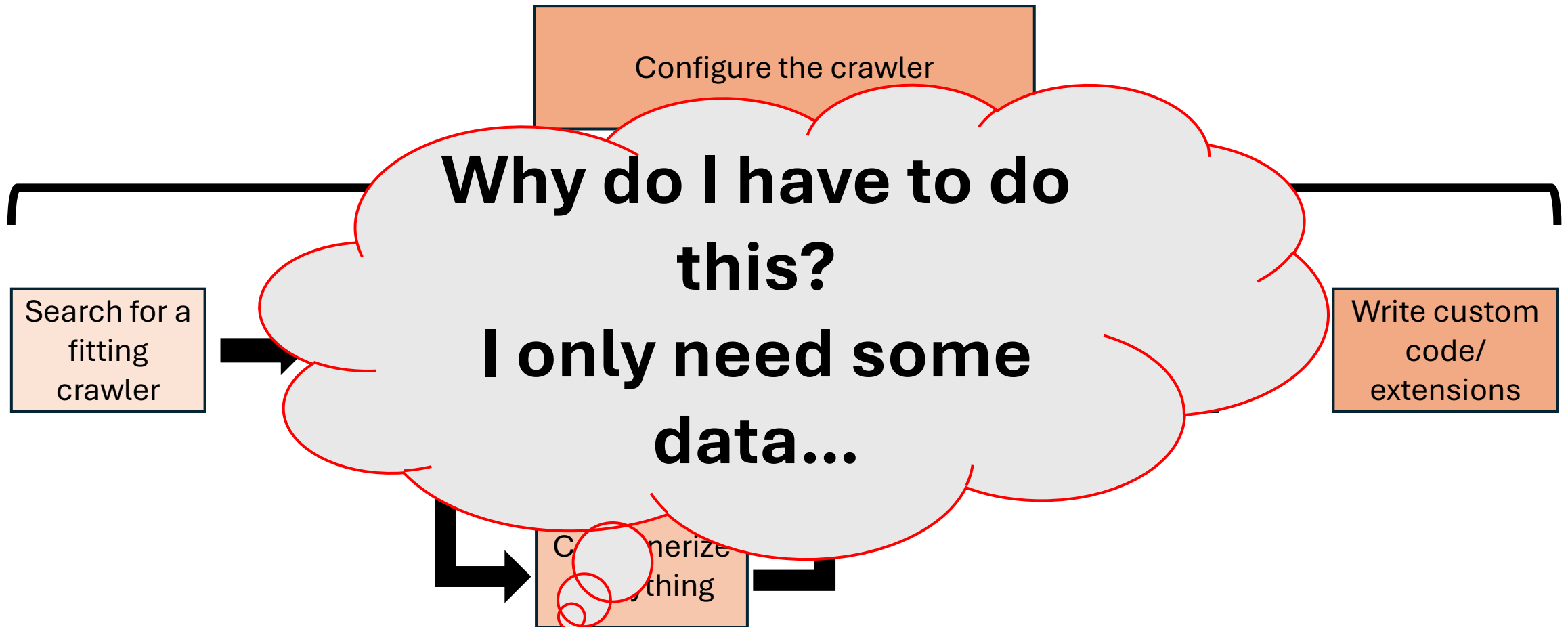


# Motivation

Configure the crawler



# Motivation



# Motivation

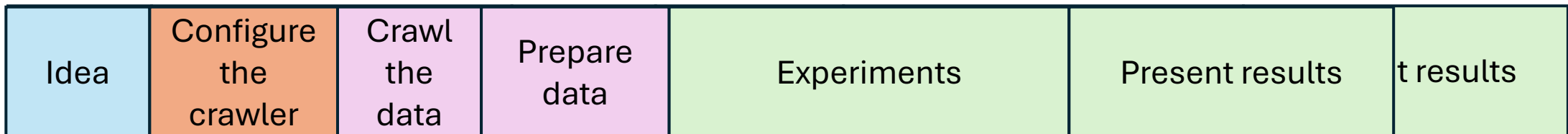
## Perfect World



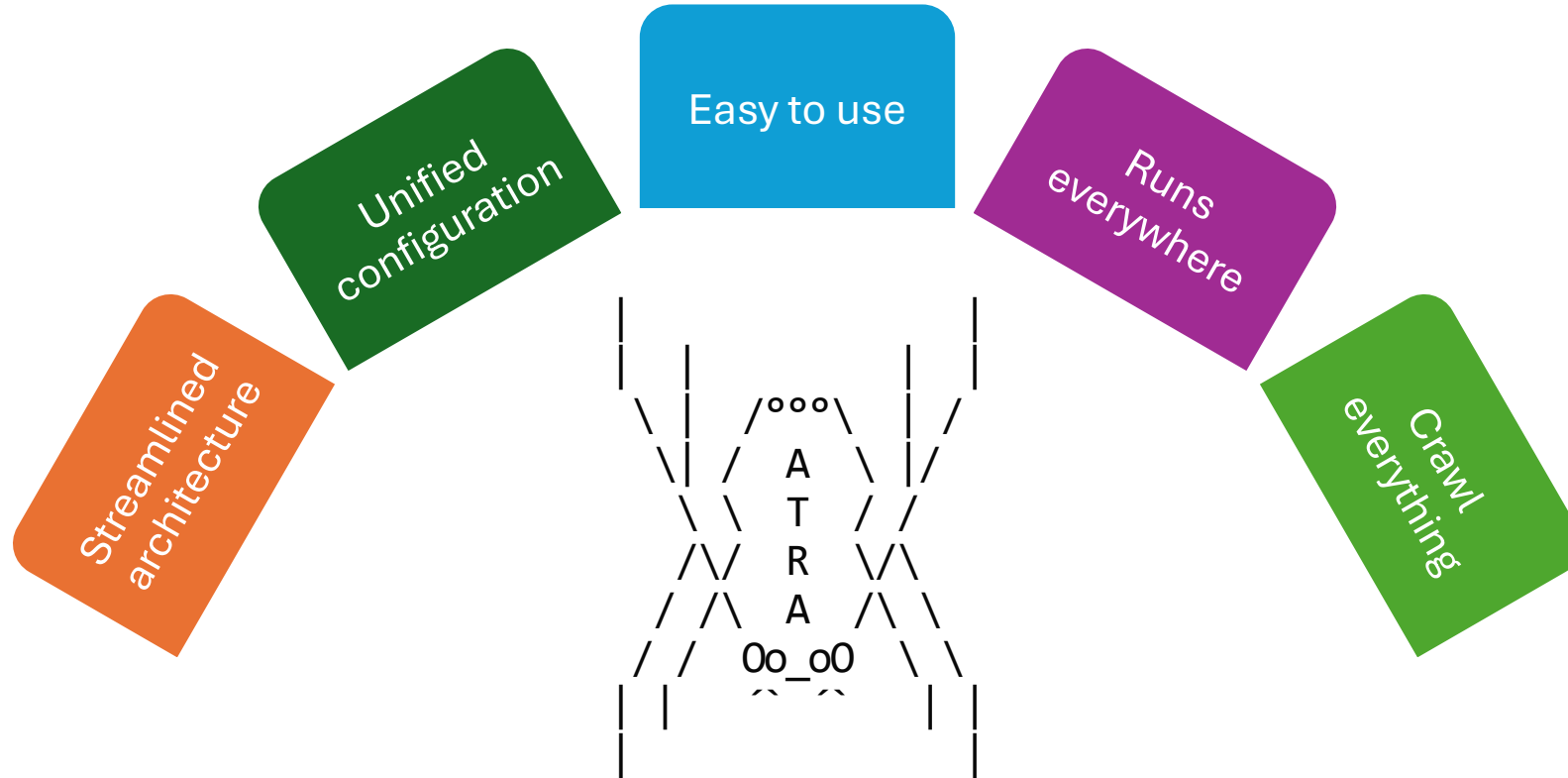
## The usual for experts



## New to crawling, but with Atra



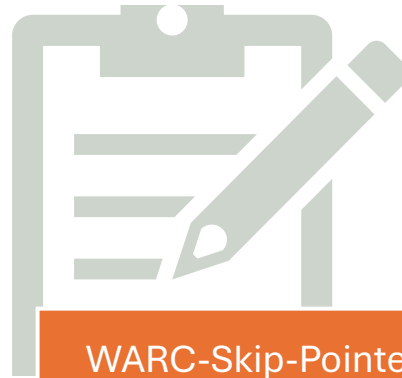
# Atra



# Key Features



Efficient Resource Usage



WARC-Skip-Pointer



Link Extraction



Encoding Support

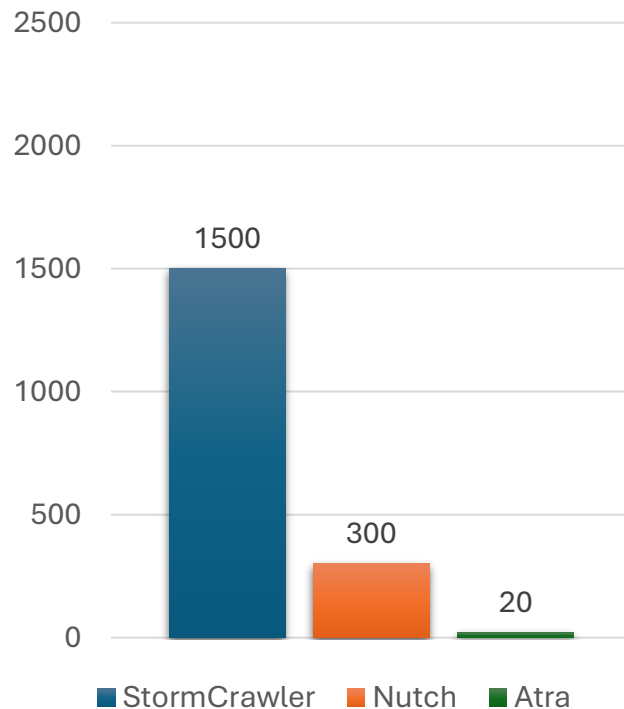


GDPR Filter

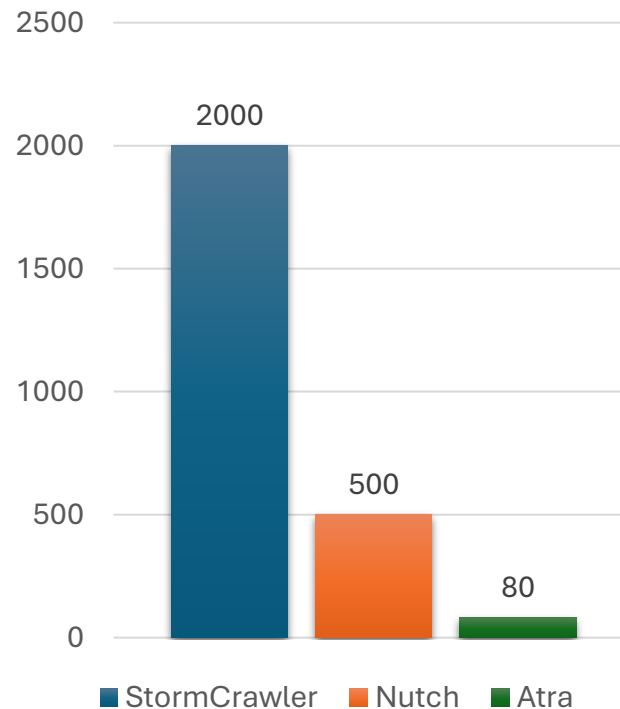


# Efficient Resource Usage

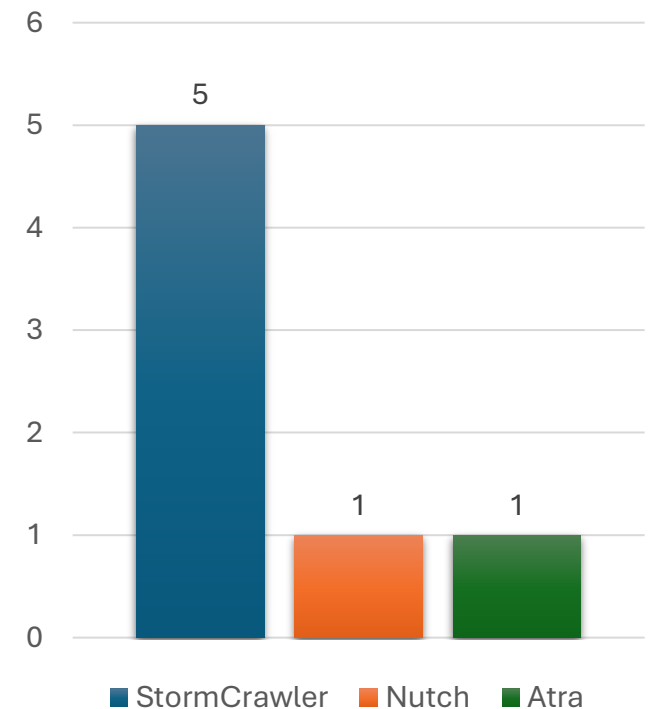
## Minimum Memory Usage (MB)



## Maximum Memory Usage (MB)



## Number of Docker Services





# WARC Skip-Pointer

WARCSkipPointer
+ path: Path + file_offset: u64 + warc_header_offset: u32 + http_header_offset: u32 + body_octet_count: u64

# Link Extraction



EXIF ODF JSON TXT OOXML  
PDF CSS XML  
SVG Binary HTML  
XLINK JavaScript RTF

# Link Extraction – HTML



## The usual way

### Apache StormCrawler

```
final Elements links = jsoupDoc.select("a[href]");
```

### Apache Nutch

<a>, <area>, <link>  
<frame>, <iframe>, <script>, <img>, <source>  
<form>

## Atra

```
/// A matcher for href locations
2 usages  ± Felix Engl +1
pub static HREF_LOCATION_MATCHER: Lazy<Regex> =
  Lazy::new(|| Regex::new(re: "location\\|s*\\.\\.s*href\\|s*=\\|s*'\\|s*([\\^']*)*\\|s*'\\|s*;?")

// Ignore [ping] of area/a
static_selectors! {
  pub [
    BASE = "base"
    HREF HOLDER = "a,area,link"
    SRC HOLDER = "audio,embed,iframe,img,input,source,track,video"
    SCRIPT HOLDER = "script"
    ON_CLICK = "[onclick]"
    FORM HOLDER = "form[action]"
    META_NO_FOLLOW = "meta[name=\\\"robots\\\"] [content=\\\"nofollow\\\"]"
  ]
}
```

# Link Extraction - Binary



- Phase 1:
  - Scan the data source and keep consecutive valid UTF-8 bytes in a **buffer**.
  - If the **buffer** fulfills one of these conditions, go to Phase 2.
    - End with "://".
    - Contains "." and has at least 4 chars.
  - If the scan encounters invalid codepoints clear the **buffer**.
  - Continue scan.
- Phase 2:
  - Consume valid UTF-8 codepoints from the source until reaching a whitespace or an invalid codepoint.
  - Use linkify to extract valid links.
  - Clear the **buffer**.
  - Go to Phase 1.

# Encoding support



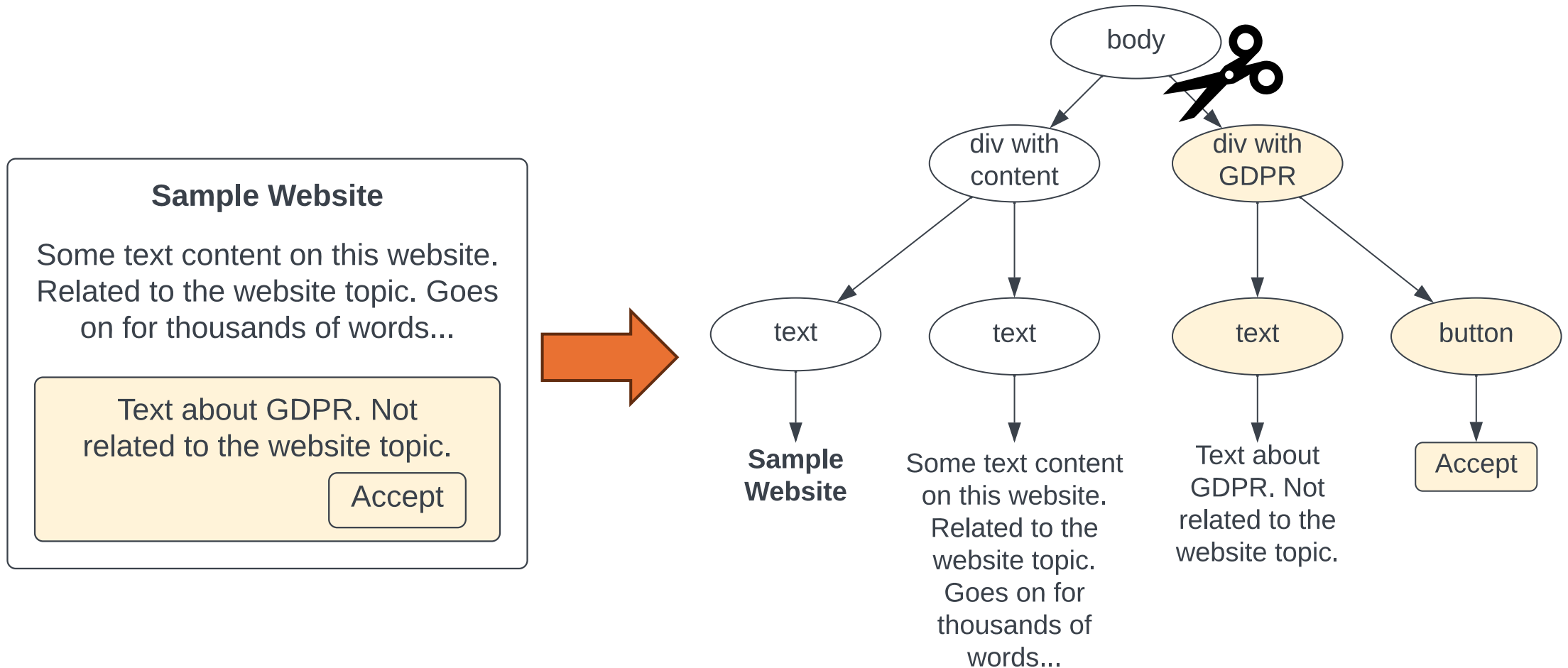
**ISO-2022-JP**  
**ISO-8859-2** **ISO-8859-5** **ISO-8859-13** **ISO-8859-8-I**  
**windows-1257** **windows-1255** **windows-1254** **KOI8-U**  
**GBK** **windows-1252** **windows-1250** **EUC-KR**  
**EUC-JP** **Big5** **Shift\_JIS** **windows-1251** **UTF-16LE**  
**IBM866** **windows-1253** **UTF-16BE**  
**windows-1258** **KOI8-R** **macintosh** **windows-1256**  
**ISO-8859-3** **ISO-8859-4** **ISO-8859-6**  
**ISO-8859-15** **ISO-8859-8** **ISO-8859-7** **x-mac-cyrillic**  
**gb18030**

# Encoding support



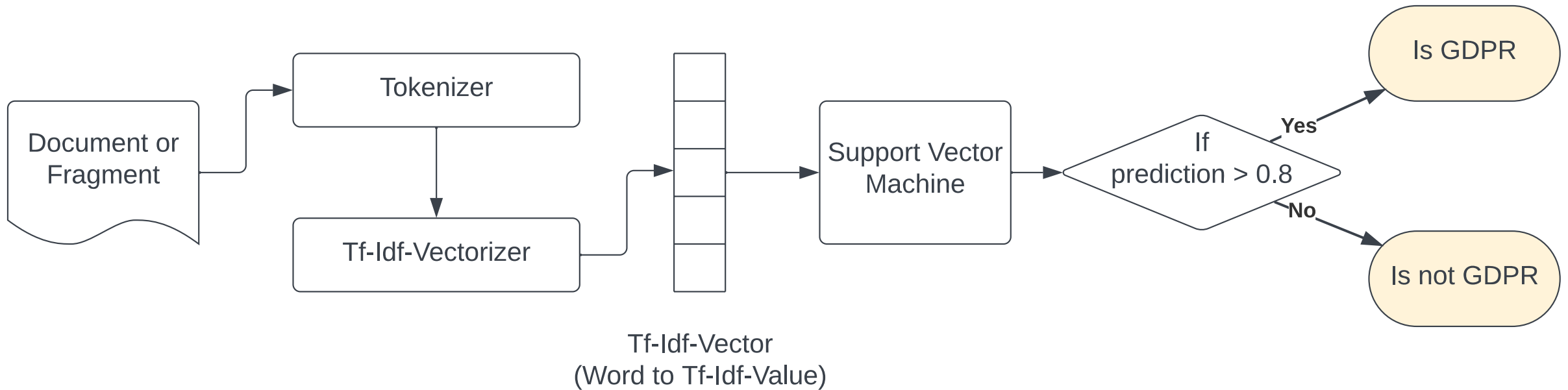
Atra uses the same decoding mechanism as Gecko Version 54+

# GDPR Filter

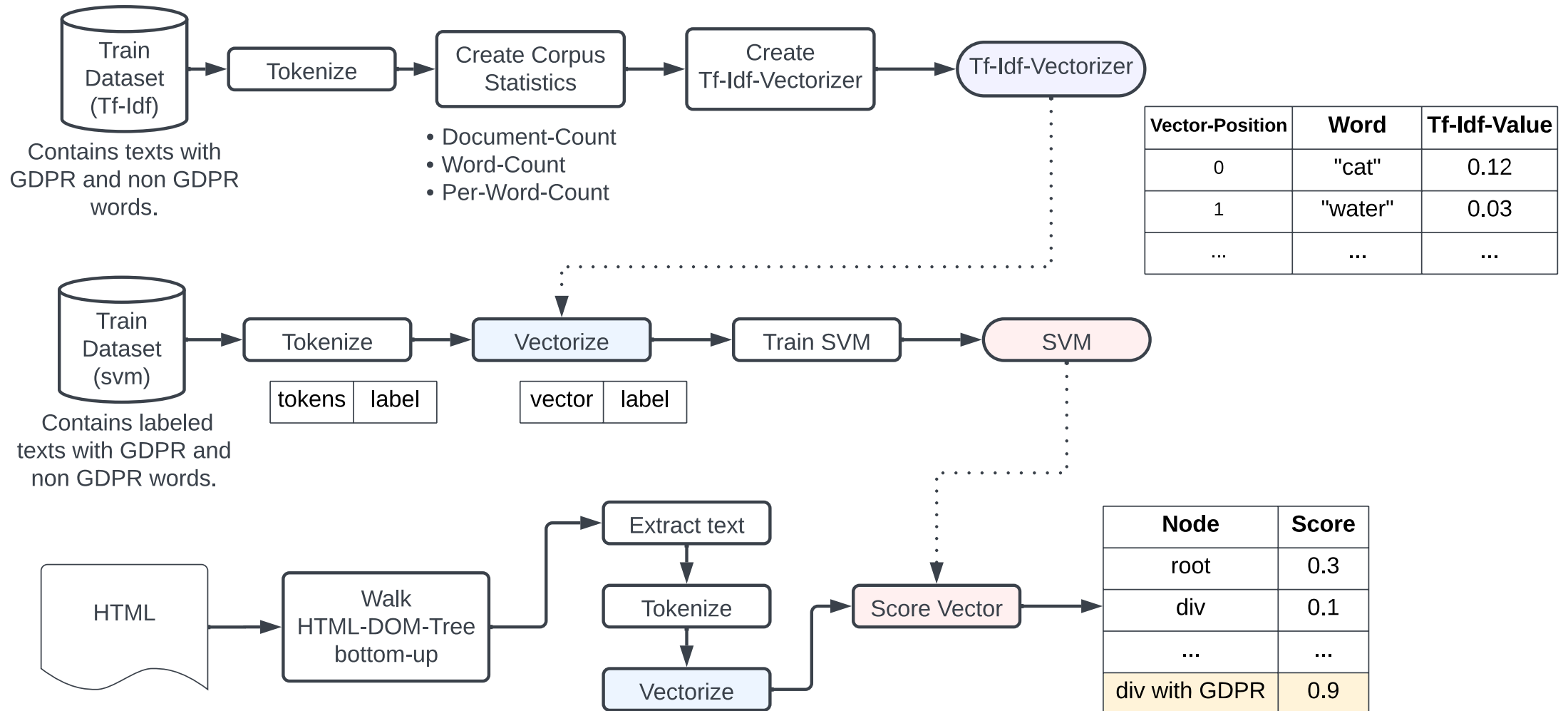




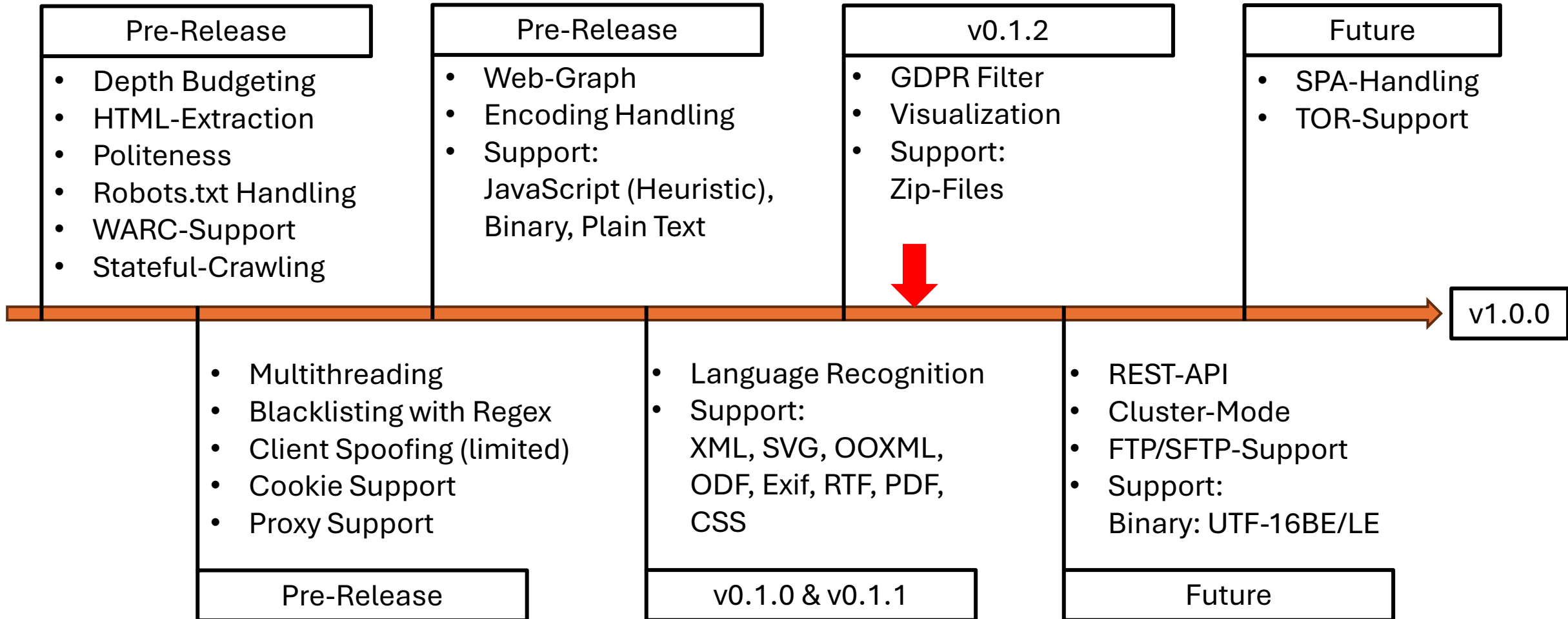
# GDPR Filter



# GDPR Filter



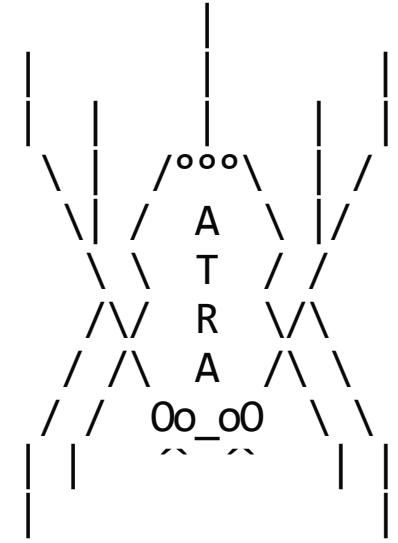
# Roadmap



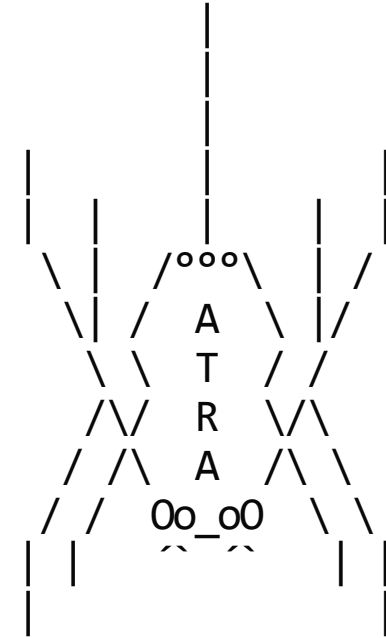
# Demo



<https://github.com/FelixEngl/atra>



Thank you for your attention!



Questions?