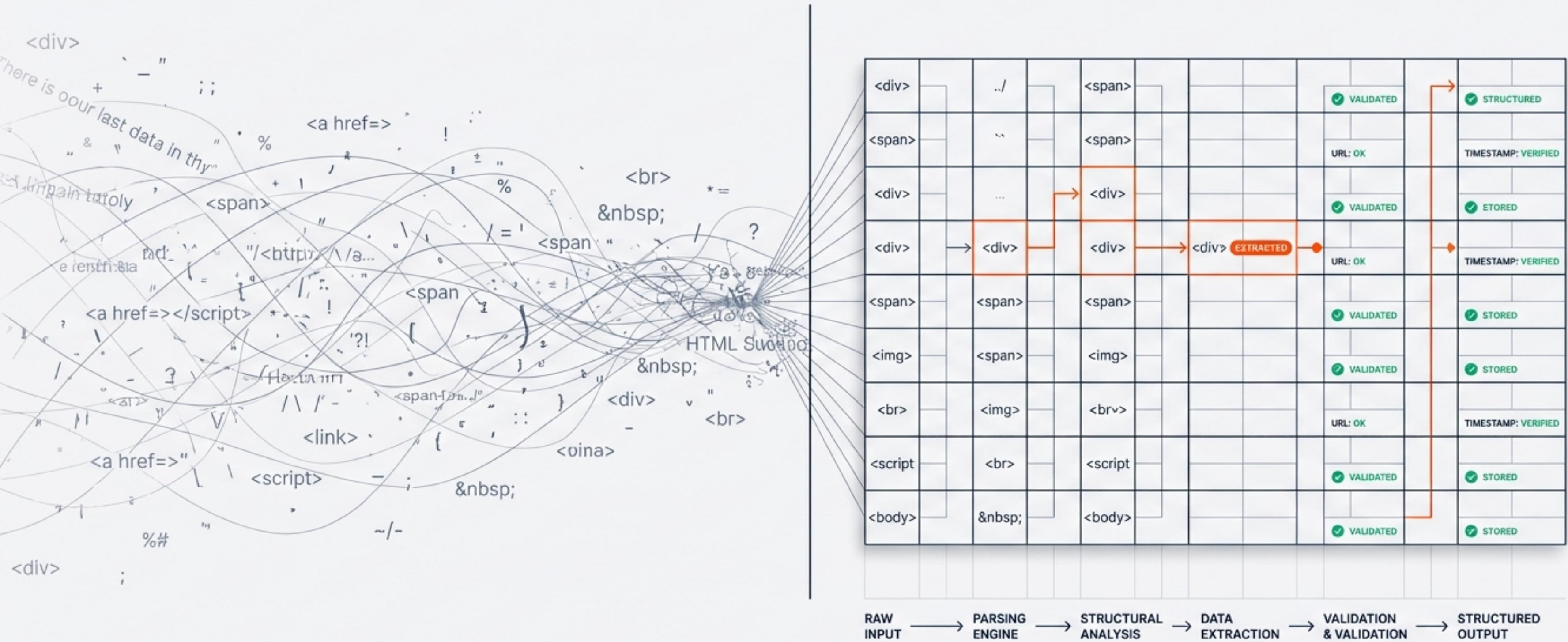


HTML Content Processing Pipeline

End-to-End Web Capture, Storage, and Intelligent Transformation



We have built a foundational engine for automated content processing at scale.

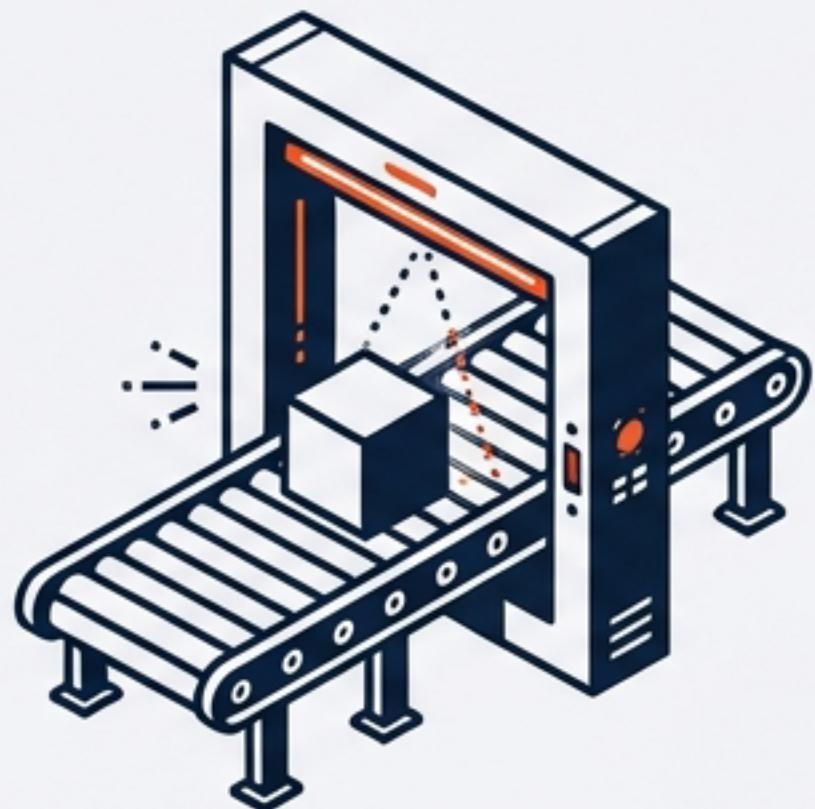
This is an end-to-end pipeline that captures web pages, secures them in a structured cache, and applies intelligent transformations to extract and filter content.



1. Capture



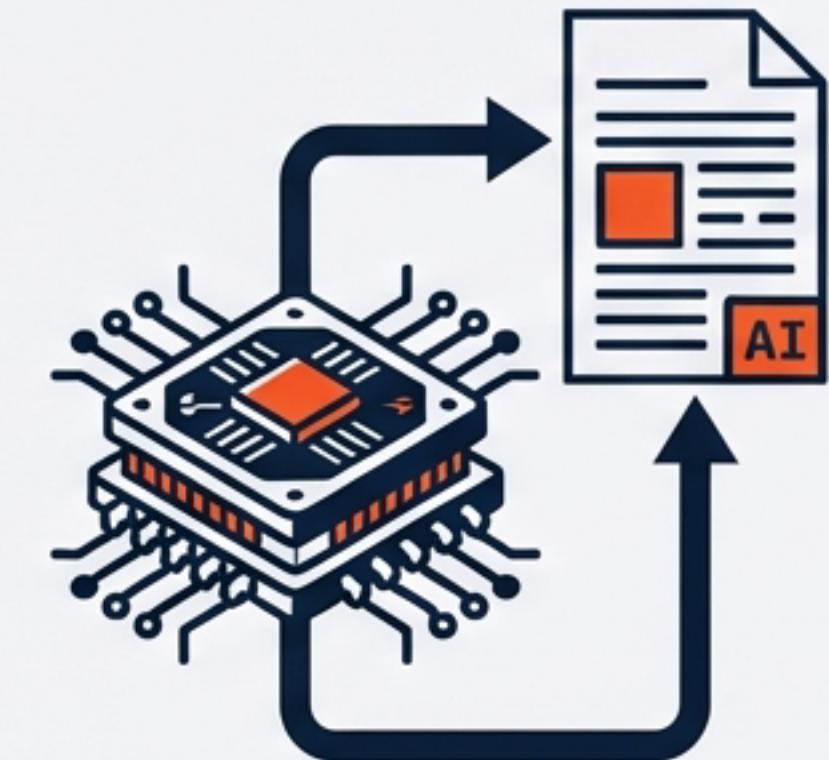
Fetching HTML from any public URL.



2. Refine



Structured storage via the LETS pipeline.



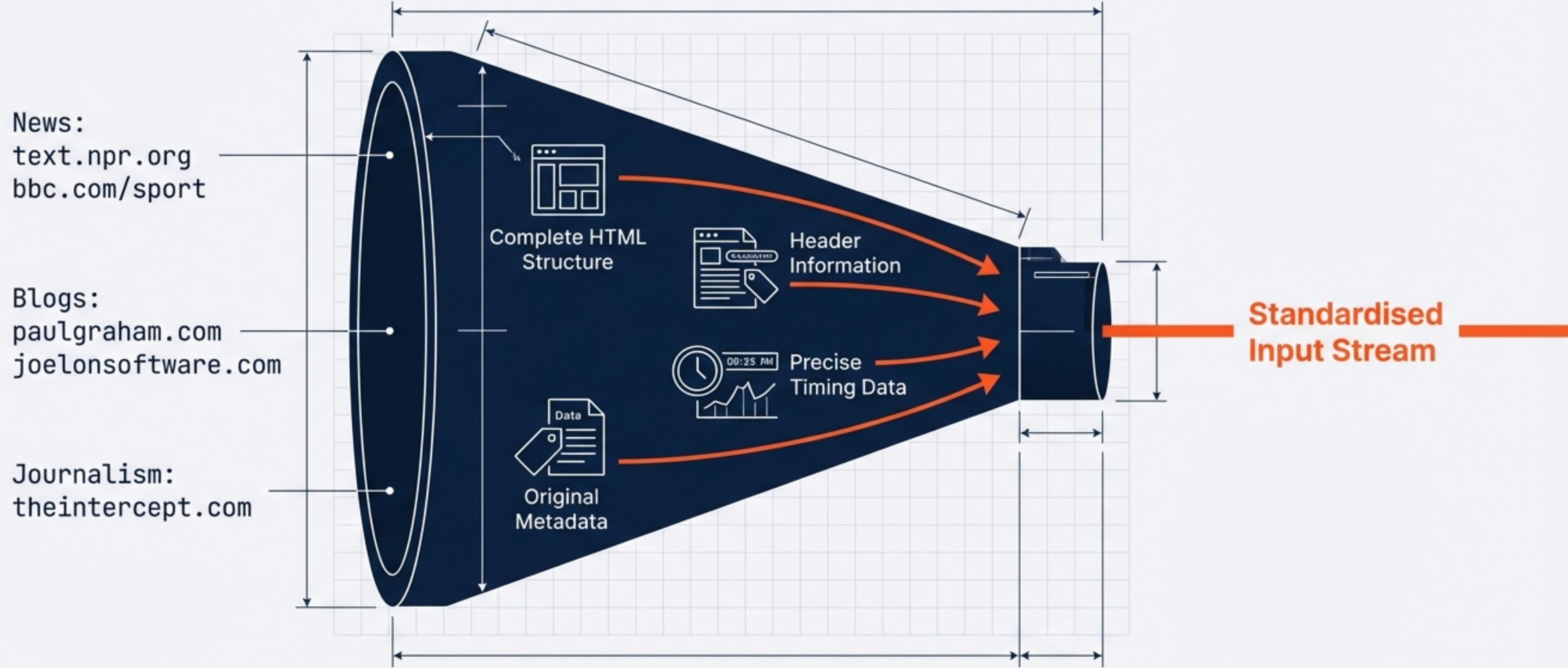
3. Utilise



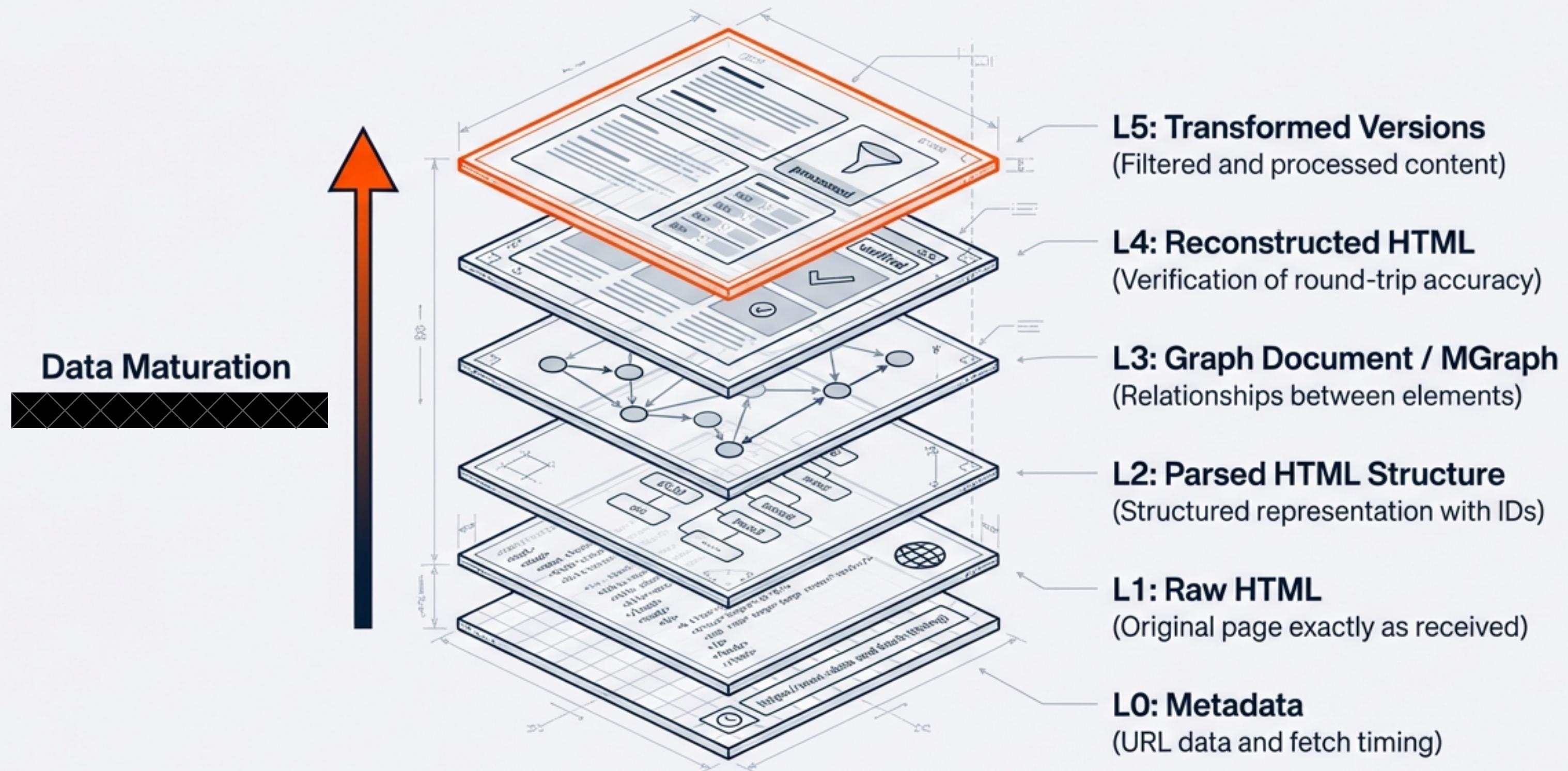
Intelligent extraction for AI and analysis.

Key Takeaway: This system replaces ad-hoc scraping with a robust architecture designed for high-fidelity data preservation and retrieval.

Ingesting diverse content sources with full context preservation

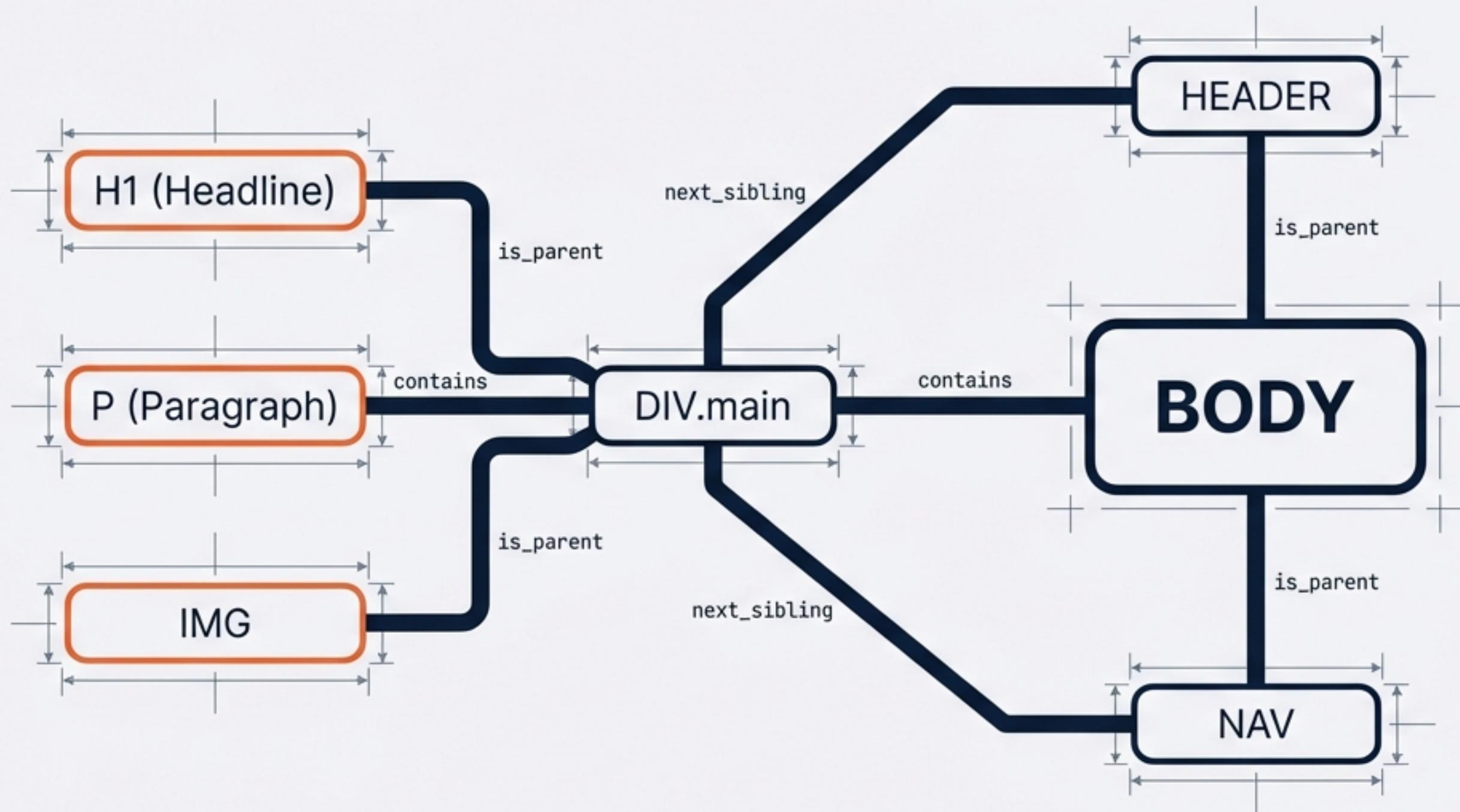


The LETS Pipeline: A multi-layer architecture for data integrity.



Key Takeaway: The LETS pipeline ensures data integrity through a structured, multi-layered maturation process from raw capture to transformed, high-fidelity outputs.

Beyond text: Understanding structure via Graph Technology.



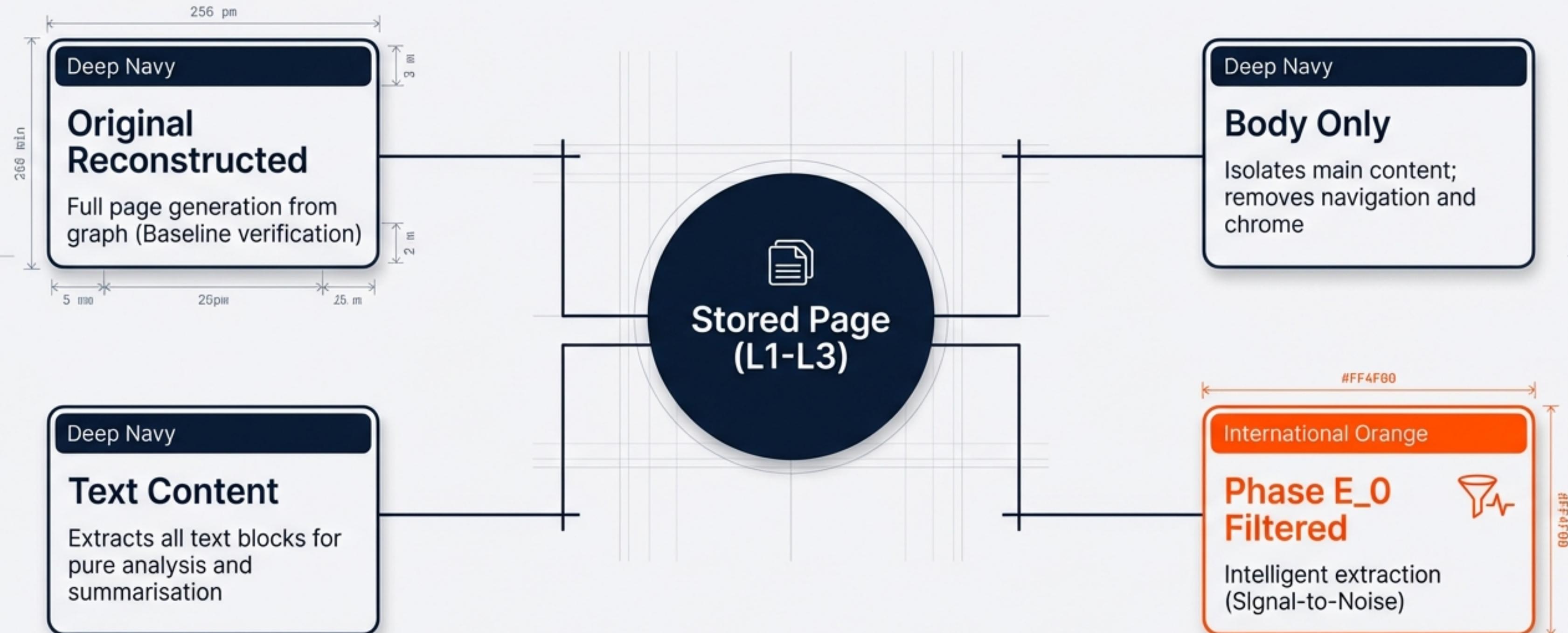
Layer L3 (MGraph) maps the relationships between every element on a page.

Why It Matters:

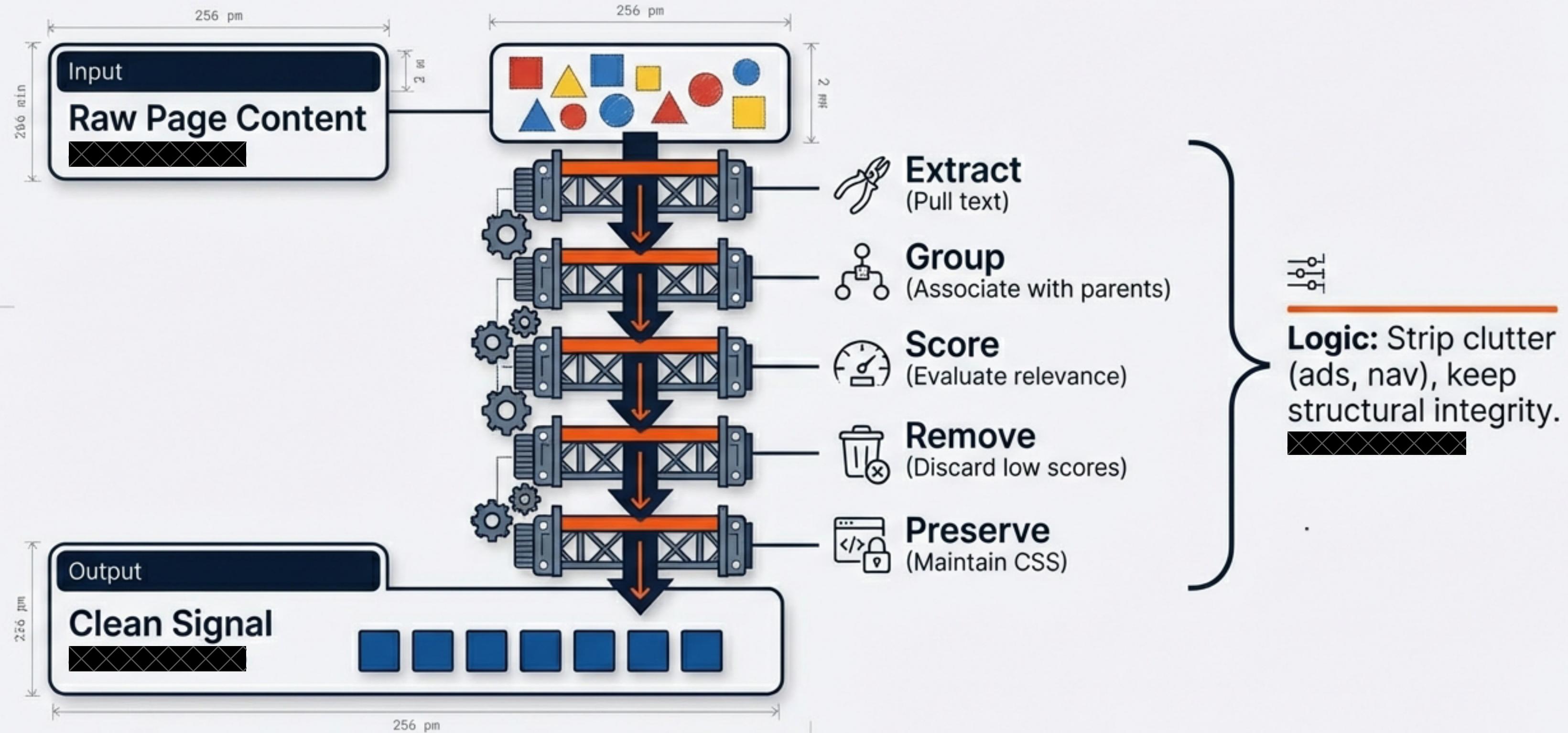
- Precise element addressing
 - Structural analysis independent of visual rendering
 - The ability to reconstruct the page dynamically

One capture, infinite views.

Once stored, we apply intelligent transformations without re-fetching from the source.



Phase E_0: The Content Filtering Engine.



Total transparency via the Cache Browser Console.

Explore, inspect, and compare cached data with precision.

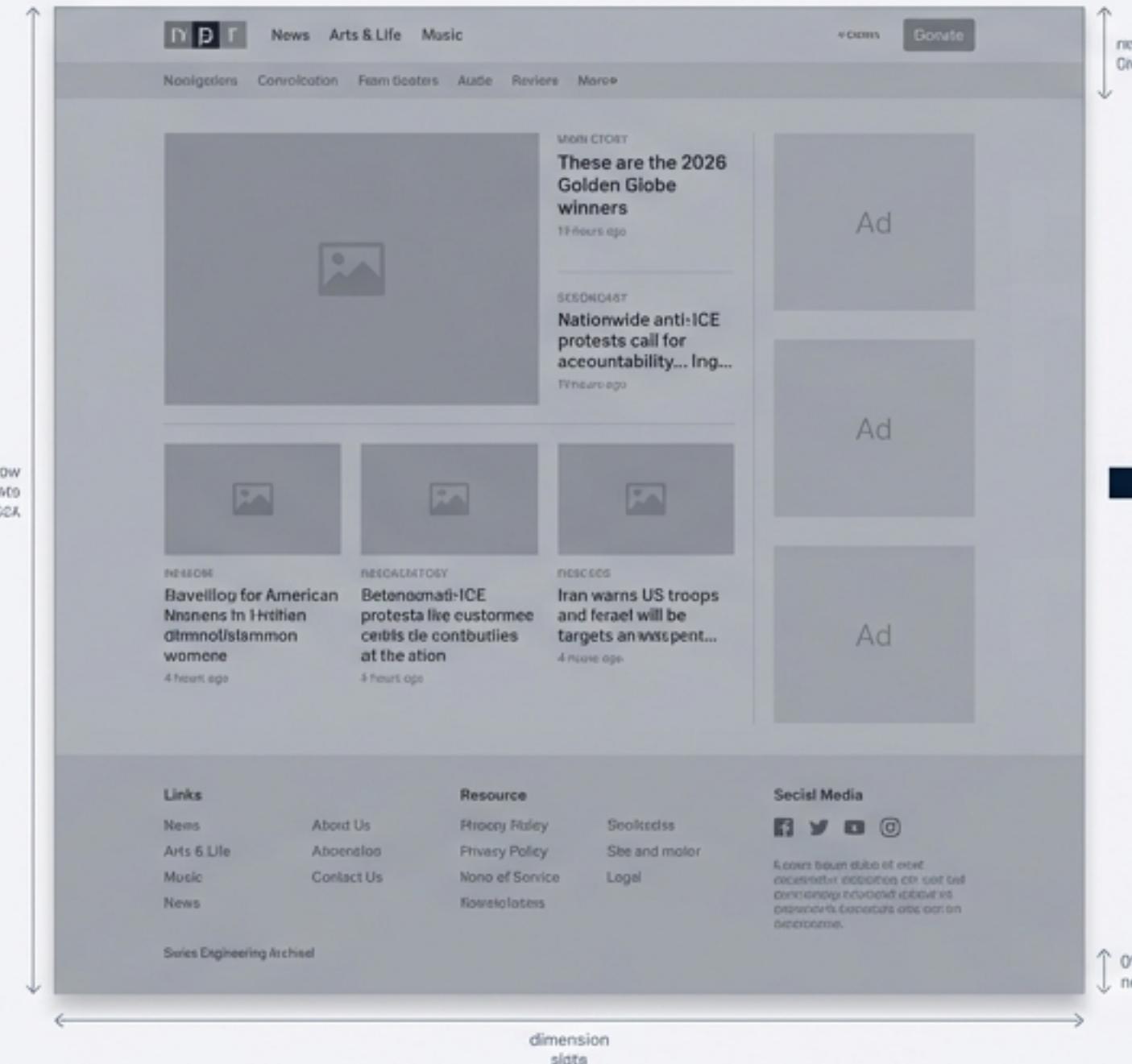
The screenshot illustrates the Cache Browser Console interface, which provides a comprehensive view of cached data across five levels of abstraction (L1 to L5). The interface is divided into several sections:

- Browse:** Navigate folder structures. The Explorer panel shows a tree view of cached data, including a file named "index.html" under the "cache_v0.1.3/mgraph_data/L5_semantic" path.
- Preview:** Render HTML directly. The Preview panel displays the "mGraph Insight" report, featuring sections like "Executive Summary" and "Key Findings", along with a bar chart titled "Data Points".
- Compare:** Toggle L1/L5 versions. A switch in the Properties panel allows users to compare different levels of cached data.
- Inspect:** View raw JSON. The Properties panel also displays the raw JSON representation of the selected data object, which includes fields such as "id", "version", "timestamp", "type", "source_url", "content", and "l1_reference".

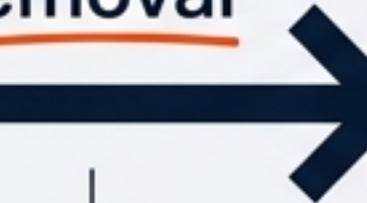
```
{  
  "id": "page_L5_7a8b9c",  
  "version": "v8.1.3",  
  "timestamp": "2024-05-21T14:38:882",  
  "type": "semantic_processed",  
  "source_url": "https://example.coa/report",  
  "content": {  
    "title": "mGraph Insight Report",  
    "summary": "Analysis shows positive trend...",  
    "data_points": [  
      {  
        "label": "Revenue",  
        "value": 1500,  
        "change": "+12%"  
      },  
      {  
        "label": "Active Users",  
        "value": 35000,  
        "change": "+5%"  
      }  
    ],  
    "l1_reference": "page_L1_1x2y2r"  
  }  
}
```

Case Study: Signal extraction on NPR.org.

Input (Raw Web)



Noise
Removal



Output (Phase E_0 Filtered)



Validated across diverse web architectures.

text.npr.org

News Headlines



paulgraham.com

Essays/Blog



theintercept.com

Investigative Journalism



joelonsoftware.com

Tech Blog



bbc.com/sport

Sports News



docs.diniscruz.ai

Technical Documentation



Business Value Delivered Today.



Automated Capture

Systematic ingestion of web content replaces manual effort.



Structured Storage

Enabling multiple views and queries on the same data.



Content Filtering

Delivers cleaner, higher-quality data for downstream use.

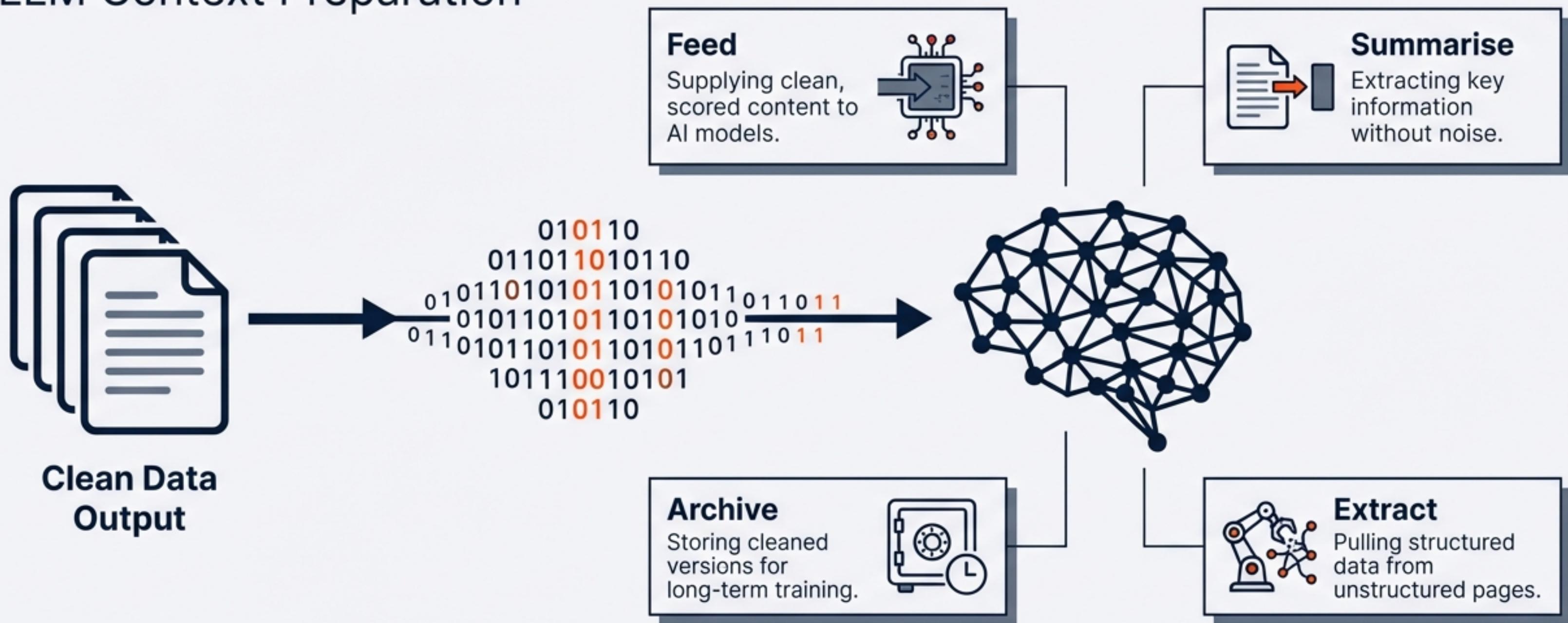


Visual Inspection

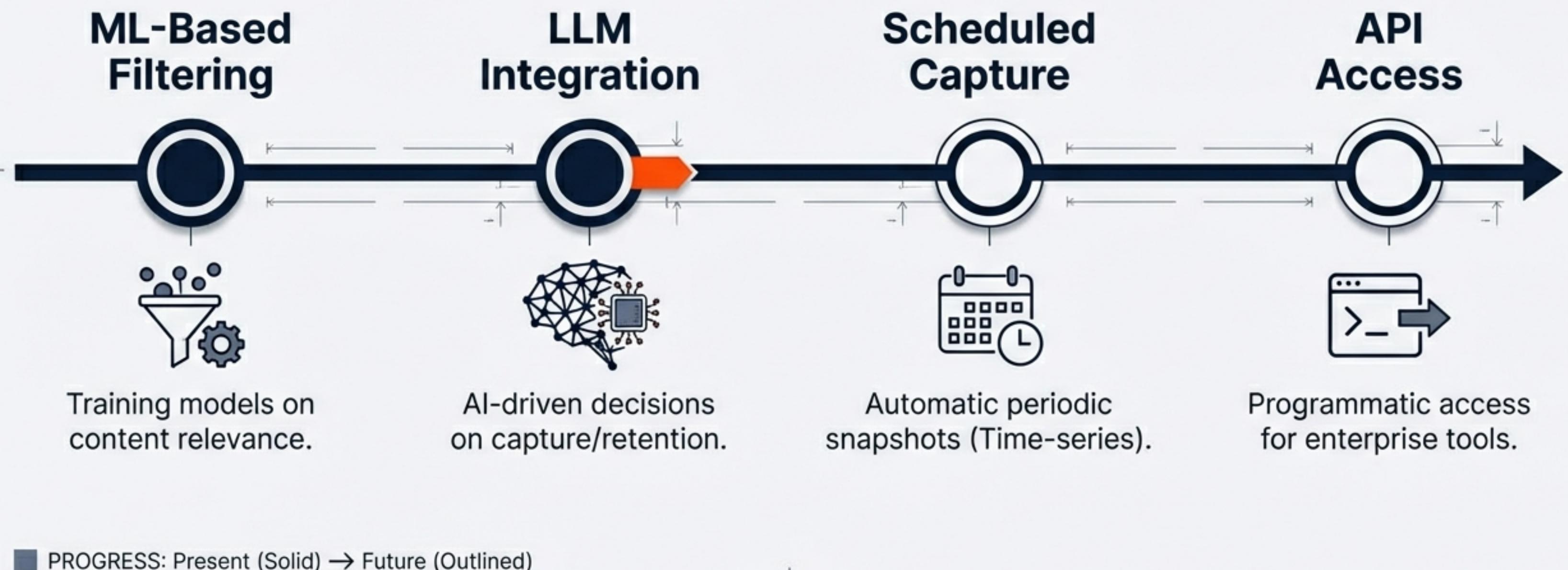
Immediate QA and verification through the Cache Browser.

The Foundation for AI Readiness.

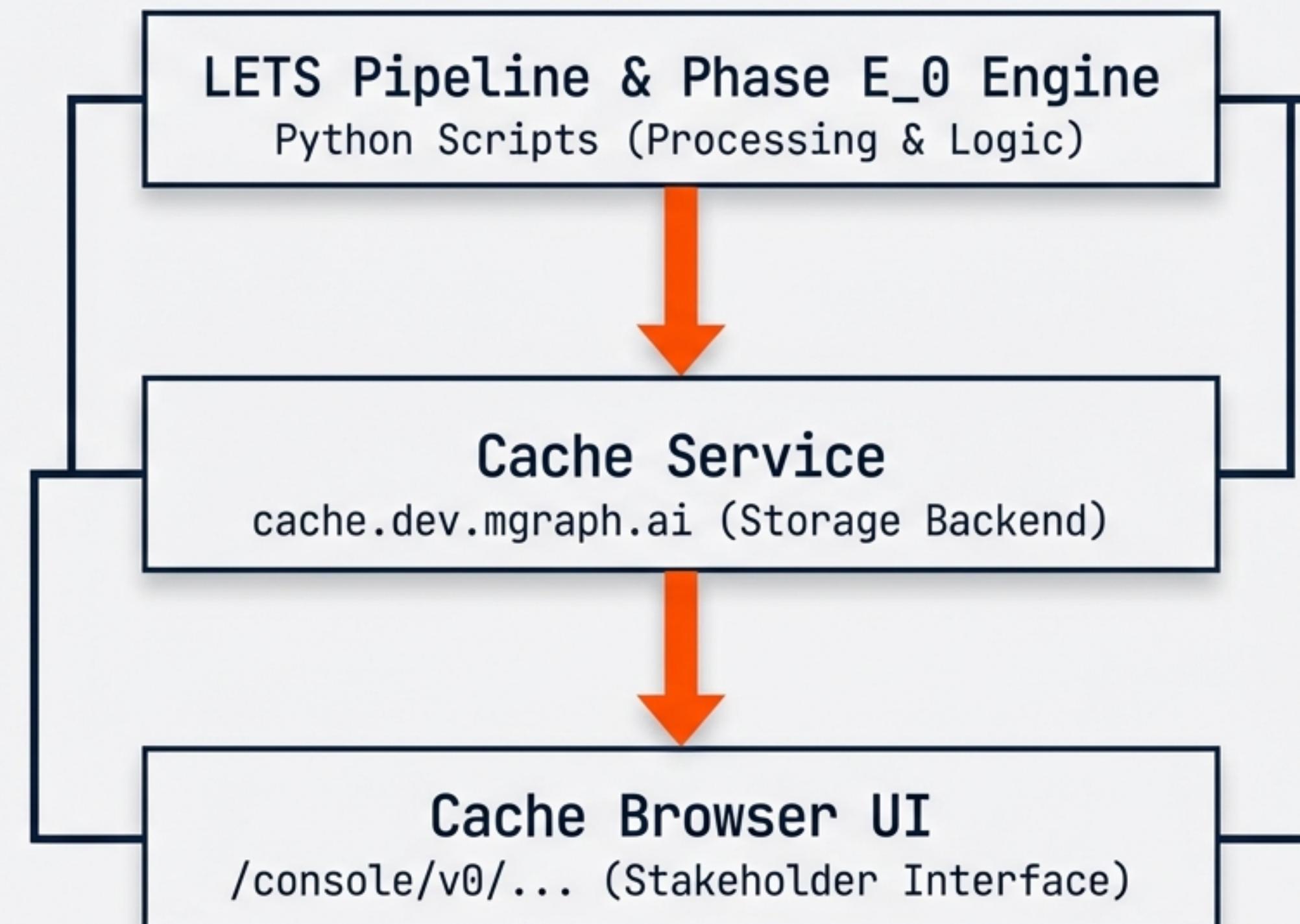
LLM Context Preparation



Roadmap: From Static Capture to Intelligent Monitoring.



Technical Infrastructure & Deployment.



Ready to Scale.

	Capture: Any public web page.
	Store: Structured, queryable cache.
	Transform: Multiple processing options.
	Filter: Intelligent Phase E_0 reduction.
	Visualise: Web-based inspection.

This infrastructure provides the foundation for automated content processing, AI-ready data preparation, and scalable web archival.