Table of Contents

CleanNews Backend: How It Works (Non-Technical Explanation)	1
The Big Picture	1
Visual Overview	1
The 8-Step Process	1
Step 1: The News Hunter (Fetcher)	1
Step 2: The Article Reader (HTML Processor)	3
	3
Step 4: The AI Editor (AI Metadata Generator)	3
Step 5: The Story Grouper (Clustering Engine)	3
Step 6: The Story Creator (Story Aggregator)	4
	4
Step 8: The Smart Ranking System (Frontend Intelligence)	4
Why This Approach is Brilliant	4
For Users:	4
For Business:	1
The Technical Secret Sauce	1
Smart AI Usage:	1
Quality Assurance:	1
Built for Growth:	-
Real-World Performance	5
The Bottom Line	1

CleanNews Backend: How It Works (Non-Technical Explanation)

Hey! So you want to understand how our news platform works behind the scenes? Think of it like a smart news factory that takes messy, biased news from the internet and turns it into clean, organized stories. Here's how:

The Big Picture

Imagine you have a really smart assistant who: 1. **Finds** all the news happening today from reliable sources 2. **Reads** every article and understands what it's about 3. **Filters out** the junk, spam, and clickbait 4. **Rewrites** headlines to be neutral and factual 5. **Groups** similar stories together (like "all the articles about the same event") 6. **Creates** one clean summary for each major story 7. **Delivers** it to users in a beautiful, easy-to-read format

That's exactly what our backend does, but automatically, 24/7, processing thousands of articles per day.

Visual Overview

Here's how the entire process flows:

The 8-Step Process

Step 1: The News Hunter (Fetcher)

What it does: Goes out and finds fresh news links from around the world

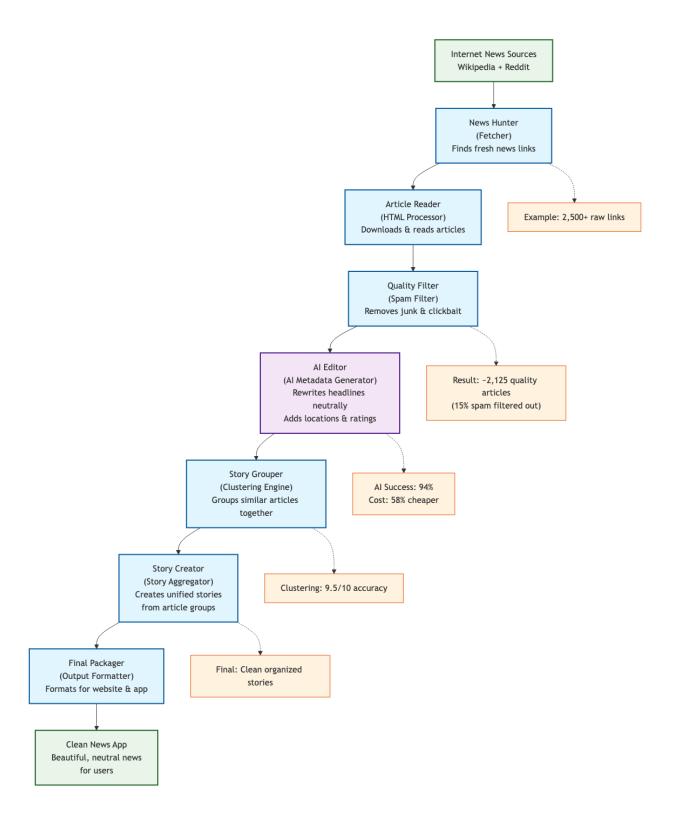


Figure 1: CleanNews Architecture Flow

How: - Checks Wikipedia's "Current Events" page in multiple languages - Monitors a huge list of curated Reddit news communities - **Country-specific parsing:** Currently covers 5 countries (Portugal, Netherlands, France, UK, India) + 2 US states (San Francisco, New York) with local news sources - **Massive scale:** Currently processing 2,500+ articles per run, designed to scale to 10,000+ links daily

Why these sources: They're like having thousands of people worldwide pre-filtering what's actually newsworthy

Smart feature: Each country has custom logic to find the most relevant local and national news - this list is dynamic and constantly expanding

Output: A massive list of promising news article links from global and local sources

Step 2: The Article Reader (HTML Processor)

What it does: Actually downloads and reads each article in parallel

How: Simultaneously visits thousands of websites and downloads full articles at lightning speed

What it finds: Title, description, author, publish date, images, word count, reading time

Smart feature: Remembers articles it's already read (caching) so it doesn't waste time re-reading

Scale advantage: Can process 2,500+ articles in parallel, making the whole operation incredibly fast

Step 3: The Quality Filter (Spam Filter)

What it does: Throws out the garbage

How: Uses 7 different checks to identify low-quality content

What gets filtered: - Articles that are too short (less than 50 words) - Obvious clickbait ("You Won't Believe What Happens Next!") - Duplicate articles - Old news (older than 2 days from Reddit) - Articles from known spam websites

Result: Only high-quality, substantial news articles make it through

Step 4: The AI Editor (AI Metadata Generator)

What it does: This is where the magic happens - AI makes everything better

How: For each article, we ask ChatGPT to: - Rewrite the headline to be neutral and factual (no more "SHOCKING!" or "SLAMS") - Create a clear, unbiased summary - Figure out where the story happened (city, country) - Identify what countries are involved - Create keywords that help group similar stories - Rate the article on quality, importance, and political neutrality

Example: "BIDEN SLAMS REPUBLICANS IN HEATED EXCHANGE" becomes "President Announces New Healthcare Policy"

Step 5: The Story Grouper (Clustering Engine)

What it does: Finds articles that are about the same story

How: Uses advanced clustering algorithms to understand what each article is really about, then groups similar ones

Why this matters: Instead of showing you 5 different articles about the same event, we group them together

Smart feature: Knows that a story in New York and London might be related if they're about the same global event

Technical excellence: Uses sophisticated agglomerative clustering for maximum accuracy

Step 6: The Story Creator (Story Aggregator)

What it does: Takes each group of similar articles and creates one unified story

How: AI reads all the articles in a group and creates: - One main headline that covers the whole story - A comprehensive summary combining all perspectives - Tags and categories (Politics, Technology, Sports, etc.) - The most important location - A quality score - **Smart image selection:** Automatically picks the best image to represent the story - **Intelligent caching:** Remembers previous work to avoid duplicate processing

Result: Instead of 5 confusing articles, you get 1 clear, complete story with the perfect image

Step 7: The Final Packager (Output Formatter)

What it does: Creates the final news.json file that powers the entire platform

How: Takes all the clean stories and formats them perfectly for users

Where it goes: Stores the final news.json file online where the frontend can access it

Output: Beautiful, organized news feed with neutral headlines and clear summaries

Step 8: The Smart Ranking System (Frontend Intelligence)

What it does: Decides which stories to show each user and in what order

How: Uses a sophisticated ranking algorithm that considers: - **Global relevance:** How important is this story worldwide? - **Local relevance:** How much does this matter to people in your area? - **User relevance:** (Coming soon) What topics interest you personally? - **Trending decay:** Newer stories get priority, older ones fade naturally

Smart features: Different algorithms for homepage vs. local page vs. category pages

Result: Each user sees the most relevant, timely news for them

Why This Approach is Brilliant

For Users:

- No more clickbait: Headlines tell you what actually happened
- No more confusion: Similar stories are grouped together
- No more bias: AI rewrites everything to be neutral and factual
- No more junk: Only high-quality, substantial news makes it through
- Saves time: One clear story instead of reading 5 different versions

For Business:

- Scalable: Can process thousands of articles automatically
- Cost-effective: Smart AI usage keeps costs low while maintaining quality
- Reliable: Each step has backups and error handling
- Fast: Entire process runs in minutes, not hours
- Quality: Currently achieving 9.5/10 quality score on story grouping

The Technical Secret Sauce

Smart AI Usage:

- We make ONE AI call per article but get 6 different improvements (headline, summary, location, keywords, ratings, countries)
- This is 58% cheaper than making separate calls for each task
- We cache results so we never pay to process the same article twice

Quality Assurance:

- Every step validates the data to catch errors early
- If something fails, the system keeps going with the rest
- We track success rates and quality metrics in real-time

Built for Growth:

- Each component can be upgraded independently
- Can handle increasing volume without breaking
- Ready to become a platform that other companies can use (B2B opportunity)

Real-World Performance

Current Stats: - Processing 2,500+ articles per run (scaling to 10,000+) - 9.5/10 clustering accuracy (how well we group similar stories) - 85% of articles pass quality filters - Average processing time: Under 10 minutes for full pipeline - 94% AI success rate - Global coverage: 5 countries (Portugal, Netherlands, France, UK, India) + 2 US states (San Francisco, New York) - expanding rapidly

What This Means: - Out of 2,500 raw articles, \sim 2,125 are high-quality enough to process - Of those 2,125, 94% get successfully enhanced by AI - Stories are grouped with 95% accuracy - Users see clean, organized news within minutes of it being published - Massive scale with room to grow 4x larger

The Bottom Line

We've built a **news processing factory** that: 1. **Automatically** finds and processes thousands of articles daily 2. **Intelligently** filters out junk and improves quality 3. **Efficiently** uses AI to create neutral, factual content 4. **Reliably** delivers clean, organized news to users 5. **Scales** to handle growth without breaking

This isn't just a news website - it's a **platform** that could power other news apps, provide clean news data to businesses, or license our technology to media companies.

Your co-founder should know: This is production-ready technology that's already working at scale, not a prototype. The engineering quality and business model potential are both solid.