

# System Design: Web Spam Classification Engine

By Isaac Bernat | [isaacbernat.com](https://isaacbernat.com) | August 2025

## Executive Summary

This document presents a technical design for a proof-of-concept (PoC) web spam classifier, developed to solve a common real-world challenge under two key constraints: a legacy ecosystem (**Perl**) and a tight implementation timeframe (**20 hours**).

The proposed solution is a “white box” **heuristic engine**, chosen for its interpretability, tunability and rapid validation capabilities over a more complex Machine Learning (ML) model. The design is modular, featuring a decoupled feature extractor and a rules engine driven by a human-readable YAML configuration.

The document details a success framework which prioritizes **high precision**, a comprehensive **multi-layered testing** strategy (including adversarial testing) and a **phased path to production** that evolves the system **into** a modern **Python/ML service**. The design demonstrates a philosophy of pragmatism and strategic architectural trade-offs in a constrained environment.

## Introduction: A Pragmatic Design for a Real-World Scenario

Designing a system to accurately identify web spam is a classic backend challenge. This document presents a technical design for a PoC spam classifier. The goal was to solve a common, real-world engineering problem: **how to build a new, critical feature within the constraints of an established, legacy ecosystem**.

To simulate this scenario, the design was developed under a specific set of constraints that heavily influenced the architectural decisions:

- **A Legacy Tech Stack:** The design assumes the target environment is a large, mature **Perl monolith**. This forces a focus on integration and maintainability by a team with deep expertise in that specific stack.
- **A Tight Timeframe:** The project was scoped for a rapid PoC, with an estimated **20-hour budget for the implementation phase**. This required a design that prioritized simplicity and speed-to-delivery.
- **The Goal of Validation:** The primary objective was to validate a classification strategy against a dataset, not to build a production-ready system from day one.

These realistic constraints led to a design that deliberately prioritizes **interpretability**, **tunability** and **leveraging existing strengths** over building a “perfect” solution in a theoretical vacuum.

## Assumed Inputs

This design focuses exclusively on the classification engine. It assumes the prior existence of two key inputs:

1. **A Corpus of Scraped Websites:** A local directory containing the raw HTML content of various websites.
2. **A Labeled Dataset:** A csv file (or similar) that serves as the “ground truth,” mapping each website’s domain to a binary label (spam/not-spam).

The design also assumes this initial dataset is well-balanced and representative enough to create a useful set of initial heuristic rules. The document does not cover the implementation of the web scraper or the initial data labeling process.

## Defining Success: A Trust-Centric Approach

For a PoC, success is about proving the viability of an approach against clear measurable goals, not just delivering functional code. The core value of this classifier lies in its ability to effectively identify spam while upholding a high standard of trust.

This led to a core tuning philosophy that guided all subsequent metrics and decisions: **A False Positive** (a legitimate site incorrectly flagged) is significantly **more harmful** to user trust **than a False Negative** (a spam site that is missed). Therefore, the entire design is deliberately optimized to ensure the integrity of the results.

## Quantitative Goals

To that end, I established a clear set of success criteria:

- **Primary Metric (Precision > 90%):** Reflecting the core philosophy, the system must be highly precise. At least 90% of the sites it flags as spam must actually be spam.
- **Secondary Metric (Recall > 70%):** The system must still be effective, correctly identifying at least 70% of all true spam sites in the dataset.
- **Health Metric (Accuracy > 80%):** The overall accuracy must be significantly better than random chance, serving as a general health check.

## Tuning Levers

The design's configuration file would provide two main levers to achieve these metrics:

- **Rule Weighting:** Individual rule scores would be carefully calibrated. Rules with a higher risk of producing false positives would be given more conservative (lower) scores.
- **Threshold Setting:** The final classification threshold would be set conservatively high, directly enforcing the “precision first” principle.

## Stretch Goal: Weighted Performance

A more mature success metric, beyond the initial PoC, would be to weight the performance scores by a site's real-world impact. **Misclassifying a high-traffic, popular domain is far more damaging** than misclassifying an obscure one. A future version of the efficacy test would incorporate a public list like the Tranco Top 1M to give more weight to errors on popular sites, providing a better measure of the classifier's impact on the user experience.

This quantitative and philosophical approach provides a clear framework for evaluating the PoC and tuning its performance.

## Core Design Philosophy: Pragmatism Within Constraints

The design is guided by pragmatism, focusing on delivering maximum value and minimizing risk within the defined constraints. When considering the addition of a feature to a legacy system, stability and maintainability are paramount. This led to three core principles:

1. **Interpretability Builds Trust:** For a new, critical system like a spam classifier to be adopted by an existing team, its decisions must be explainable. I chose a “**white box**” heuristic approach where every classification can be traced back to a specific set of human-readable rules. This is crucial for gaining stakeholder trust and makes the system far easier to debug and maintain than a “black box” alternative.
2. **Tunability Enables Safe Iteration:** A system with hardcoded logic is brittle and risky to modify. My design externalizes all classification rules and their associated weights into a **simple YAML configuration file**. This decouples the core logic from the code, allowing for safe, incremental tuning and the addition of new rules without requiring a full code deployment and its associated risks.

**3. Leveraging the Existing Ecosystem:** The most critical constraint was the mature **Perl environment**. Instead of fighting this, the design embraces it as a strategic advantage. It leverages Perl's renowned text-processing capabilities, which are exceptionally well suited for a rule based engine. Furthermore, it draws architectural inspiration from the proven, time tested patterns of **Apache SpamAssassin**, a highly influential and successful open source project in the domain of heuristic filtering written in Perl. This approach ensures that the proposed solution would be easily understood, built and maintained by an existing engineering team focused in Perl, drastically reducing the project's real-world risk profile.

## The Proposed Solution: A Configurable Heuristic Engine

The proposed solution is a modular, self-contained Perl application that functions as a configurable, weighted heuristic engine. This approach directly models the methodology used by successful, time-tested systems like Apache SpamAssassin, which combine a set of distinct rules or "signals" to compute a cumulative spam score. A site is classified as spam if its total score exceeds a predefined threshold.

The application's architecture is designed for simplicity, clarity and testability, with a linear data flow and three decoupled core components.

### 1. The Feature Extractor (FeatureExtractor.pm)

This is the data collection module. It handles all I/O and parsing, with a single responsibility: take a path to a scraped website's content and return a simple dictionary (or key-value object) of well-defined features. It has no knowledge of the spam rules themselves.

The initial set of features was inspired by common spamdexing techniques and focuses on signals that are both predictive and computationally inexpensive. They fall into several categories:

- **Content-Based:** Keyword stuffing density, presence of hidden text (e.g. via CSS) and the text-to-HTML ratio to detect "thin content."
- **Link-Based:** Total count of outbound links and the ratio of external-to-internal links.
- **URL-Based:** Lexical analysis of the domain name itself, such as the count of hyphens or numbers.
- **Structural:** Presence of tags commonly used for deceptive redirects, like `<meta http-equiv="refresh">`.

## 2. The Configurable Rules Engine (RulesEngine.pm)

This is the core logic of the classifier. To provide maximum flexibility and tunability, the rules are not hardcoded. Instead, they are defined in an external, human-readable `rules.yaml` file. This completely decouples the classification logic from the application code.

A typical rule in the configuration consists of:

- **Name:** A unique, descriptive identifier (e.g. `EXCESSIVE_OUTBOUND_LINKS`).
- **Feature:** The key from the `FeatureExtractor` to inspect (e.g. `num_outbound_links`).
- **Operator:** The comparison to perform (e.g. `GREATER_THAN`).
- **Value:** The threshold to compare against (e.g. `200`).
- **Score:** The weight to add to the total spam score if the rule is triggered.

## 3. The Orchestrator (classify\_site.pl)

This is the main command line script that serves as the application's entry point. It orchestrates the entire process:

1. Parses the command line arguments (e.g. the input file path).
2. Loads the classification threshold and all rule definitions from the `rules.yaml` file.
3. Invokes the `FeatureExtractor` with the input path to get the site's features.
4. Passes the extracted features and the loaded rules to the `RulesEngine` to compute a final spam score and a list of triggered rules.
5. Compares this score against the threshold to determine the final classification. The orchestrator then produces a final result object containing the classification (spam/not-spam), the total score and the list of specific rules that were triggered, providing full explainability for the decision.

## Deep Dive: The Strategic Decision to Defer Machine Learning (ML)

The most obvious alternative to a heuristic engine is an ML classification model. While an ML approach offers high potential accuracy, I made a deliberate, strategic decision to **reject it for the initial PoC phase**. This was a pragmatic choice driven by the project's specific constraints: a tight timeframe and the context of a legacy Perl ecosystem.

The key trade-offs I considered were:

## 1. Expanding on the “Interpretability Builds Trust” Principle

As stated in the core design philosophy, a “white box” solution is crucial for building trust. The deep dive reveals why this is so critical in practice:

- **For Debugging:** A heuristic approach is fully traceable. If a site is misclassified, we can see the exact list of rules that were triggered and their scores. This makes identifying and fixing a faulty rule a simple and deterministic process.
- **For Stakeholder Buy-in:** When demonstrating the PoC to product or business stakeholders, being able to explain precisely why a site was flagged (e.g. “it has an excessive number of outbound links and hidden text”) is far more convincing than saying “the model’s prediction was 0.92”.
- **For Team Adoption:** For the existing engineering team that would inherit the system, the transparent logic of a heuristic engine dramatically lowers the maintenance burden and empowers them to contribute new rules with confidence.

## 2. Implementation Feasibility vs. Ecosystem Mismatch

The 20-hour implementation budget and the Perl technology stack made a full ML workflow infeasible and unwise.

- **A responsible ML workflow** requires a significant investment in data pipelines, feature scaling, cross-validation, hyperparameter tuning and systematic error analysis. Attempting to rush these steps would produce an untrustworthy model.
- **The Perl ecosystem**, while powerful for text processing, lacks the mature, world-class ML libraries and tooling of a language like Python. Building a robust ML pipeline in this environment would require significant foundational work, placing it well outside the scope of a rapid PoC.

## 3. Robustness in a Low-Data Environment

A PoC typically starts with a limited dataset. In this scenario, a complex ML model is at high risk of “overfitting” (learning spurious correlations that do not generalize to new websites). A carefully crafted heuristic system, based on the established first principles of web spam (as documented in [sources like Wikipedia](#) and academic research), is often more robust and performs more predictably when data is scarce.

In summary, the decision to defer ML was a conscious, pragmatic trade-off. It prioritized **speed-to-learning, interpretability and a low-risk integration path** into the existing ecosystem, the most important factors for a successful PoC.

## A Multi-Layered Testing Strategy

A classifier is only as good as the confidence we have in its results. Therefore, the design includes a comprehensive, multi-layered testing strategy to ensure correctness, robustness and performance.

**1. Unit Tests (Correctness):** These tests would validate the individual components in isolation. For example, we would provide the `FeatureExtractor` with small, self-contained HTML strings to assert that it correctly counts links and calculates ratios. We would provide the `RulesEngine` with hardcoded feature dictionaries (key-value objects) to assert that the score calculation is the expected value.

**2. Integration Tests (Cohesion):** These would test the entire application flow, from command line argument parsing to the final output. By running the main `classify_site.pl` script against temporary sample files, we could assert that the components work together as intended.

**3. Efficacy / Acceptance Testing (Performance):** This is the most critical layer to evaluate the quality of the results. The design includes a dedicated evaluation function that runs the classifier against the entire labeled training dataset. It compares the system's predictions using the "ground truth" labels and computes the final [Accuracy, Precision and Recall](#) metrics. This provides the quantitative feedback loop essential for tuning.

**4. Adversarial Testing (Robustness):** To ensure the system is resilient against real-world, messy data and adversarial attacks, the design includes a strategy for adversarial testing. This would involve a suite of tests that programmatically generate synthetic, challenging test cases:

- **Spam Injection:** This test would take known "not spam" sites from the training set and automatically inject spam signals (e.g. adding 100 instances of a keyword in a `<div>` with `style='display:none;'`). The test would then assert that the site's spam score increases as expected.
- **Evasion Testing:** This test would take known "spam" sites and attempt to break the classifier. For example, it could programmatically mutate the HTML by removing random closing tags or adding malformed attributes, asserting that the application does not crash and still makes the correct classification.

## The Path to Production: From Prototype to Integrated Service

The PoC is designed as a standalone command line tool for safe, offline validation. The journey to a robust, production-grade system involves a pragmatic, phased integration into the existing legacy monolith and a strategic evolution of the technology stack.

**Phase 1: Asynchronous Integration & Performance Optimization:** The first step is to wrap the classifier's logic into an **asynchronous worker service**. A separate crawling system would be responsible for fetching website content and storing it. The crawler (or another service) would then place a **message on a queue** containing a pointer to this stored content. The worker would consume these messages, fetch the content, **perform the classification and store the result**. This is a low risk integration pattern that prevents the classifier from adding any latency to user facing requests.

**Phase 2: Building a Dynamic Data Pipeline:** The static training set would be replaced by a dynamic data pipeline that continuously scrapes and refreshes website data. To allocate resources efficiently, this pipeline would **prioritize refreshing popular domains** with high traffic more frequently than less common ones, ensuring the classifier's data is freshest where it matters most.

**Phase 3: Strategic Evolution to a Python/ML Service:** Once the heuristic engine is proven and a rich dataset has been collected, the next logical step is to introduce Machine Learning. At this point, it becomes a sound architectural decision to build a **new, dedicated V2 service in Python** to leverage its world-class ML ecosystem (`scikit-learn`, `pandas`, etc.).

- This new Python service would initially consume the same features from the `FeatureExtractor`, allowing for a direct ("apples to apples") comparison between the heuristic model and the new ML model.
- This phased approach allows the organization to strategically introduce a new technology (Python) for a specific, high-value problem, rather than attempting a risky, large-scale rewrite.

### Production Optimizations

As the system moves to a high-volume production environment, several optimizations would be implemented:

- **Early Exit (Short-Circuiting):** After the initial rule set is validated, the engine would be modified to stop processing rules as soon as a site's cumulative score exceeds the threshold. During the PoC phase, this is intentionally disabled to ensure we collect data on **all** triggered rules for better analysis.



- **Concurrency Model:** The design is inherently efficient, as the classification process is primarily **I/O-bound** (reading files) rather than **CPU-bound** (the heuristic checks are computationally cheap). This means a single worker instance could be made more efficient by using an event-driven, non-blocking I/O model to process multiple sites concurrently, maximizing throughput.

## Evolving the Product Capabilities

Beyond the technical migration, the architecture allows for significant enhancements over time:

- **Richer Heuristics:** The modular `FeatureExtractor` is designed for extension. While the PoC focuses on simple, static HTML analysis, a production version would be enhanced to extract more sophisticated signals, such as detecting obfuscated text or analyzing the behavior of embedded scripts, to combat more advanced spam techniques.
- **Automated Rule Tuning:** With a dynamic data pipeline in place, we could implement a workflow for automatically tuning the rule weights (or model re-training) based on the latest data, ensuring the classifier continuously adapts to new spammer tactics.
- **Reputation History:** The system could begin tracking a domain's classification over time. A site that frequently flips between spam/not-spam could incur a reputation penalty, making it harder for adversarial actors to game the system.
- **Finer-Grained Classification:** While the PoC may operate on a pre-aggregated "site" level, the underlying `FeatureExtractor` works on a page-by-page basis. A major evolution would be to expose this page-level classification to handle large platforms with user-generated content (e.g. `github.io`). This would allow the system to flag a single malicious page without penalizing the entire legitimate domain.

This evolutionary path, encompassing both the technical stack and the product's capabilities, allows for a gradual, de-risked migration from a simple Perl prototype to a sophisticated, Python ML service. It ensures that value is delivered and validated at every stage, prioritizing long term stability and effectiveness.