# Prompt Injection Detection Pipeline Architecture

*Input-Side Detection Before LLM Processing*

```
INPUT          Normalizer        Signature           Fusion           DECISION      ALLOWED    LLM
Query          Unicode Fix       Detector (v1)       OR Logic         Attack/Benign            Processing
                                 Heuristic           (v1+v3)
                                 Detector (v2)                                       BLOCKED
                                 Semantic
                                 Detector (v3)
```

## Example: Attack Processing Flow

```
Input:               Normalizer:        v1 Signature:        Decision:
"Ignore previous     Cleans unicode     Matches keyword      BLOCKED
instructions..."
```

*Performance: 87% TPR on known attacks | 0.77% FAR on benign | <1ms latency (GPU)*

**Production Configuration: Normalizer + v3**

**Component Specifications**

True Positive Rate (TPR): 87% on known attacks
False Alarm Rate (FAR): 0.77% on obfuscated benign
Latency: <1ms per sample (GPU-accelerated)
Complexity: ~1,200 lines of code
Deployment: Stateless, parallelizable
Dependencies: sentence-transformers, torch

Signature Detector (v1): Keyword-matching
  • 80% TPR, 0% FAR on Phase 1 attacks
  • Catches: plain, delimiter, role confusion

Semantic Detector (v3): Pattern analysis
  • 57% TPR, 0% FAR on Phase 1 attacks
  • Catches: formatting, semantic anomalies

Fusion (OR Logic): v1 OR v3
  • Combined: 87% TPR, 0% FAR

**Key Design Principles**

1. INPUT-SIDE DETECTION: Attacks blocked BEFORE reaching the LLM, preventing prompt injection at the source
2. NORMALIZER FIRST: Unicode/homoglyph normalization ensures consistent detection across obfuscation techniques
3. COMPLEMENTARY DETECTORS: v1 (signature) + v3 (semantic) catch different attack patterns through OR fusion
4. THRESHOLD-INVARIANT: Binary OR logic eliminates threshold tuning complexity in deployment
5. PRODUCTION-READY: <1ms latency with GPU acceleration, stateless architecture

Legend:
☐ Input    ☐ Normalizer    ☐ Detector    ☐ Fusion    ☐ Decision    ☐ LLM