

# Design and Implementation of a Custom Network Proxy Server

*Async HTTP/HTTPS Forward Proxy with SNI Inspection,  
Policy Engine, and Production-Grade Hardening*

Aditya Pratap Singh

December 19, 2025

## Abstract

This document presents a complete design and implementation report for a production-style asynchronous HTTP/HTTPS forward proxy implemented in Python. The proxy supports HTTP and CONNECT/HTTPS tunneling, TLS SNI inspection (no MITM), a policy engine (domain/method/rate limiting), robust backpressure handling, timeouts and hardening against attacks (Slowloris, header bombs), layered YAML configuration, CLI, metrics endpoint, graceful shutdown, and packaging for distribution.

The report includes enhanced architecture diagrams, per-file technical breakdowns, code excerpts, and validation strategies for high-concurrency environments.

## Contents

---

<b>1</b>	<b>Problem Statement</b>	<b>3</b>
<b>2</b>	<b>High-level Design Summary</b>	<b>3</b>
<b>3</b>	<b>Architecture Diagrams</b>	<b>3</b>
3.1	Component Diagram . . . . .	3
3.2	Sequence Diagram: HTTPS CONNECT (SNI Path) . . . . .	4
<b>4</b>	<b>Implementation Overview and Repository Layout</b>	<b>5</b>
<b>5</b>	<b>Per-File Detailed Explanation</b>	<b>5</b>
5.1	run_proxy.py . . . . .	5
5.2	src/config.py . . . . .	5
5.3	src/http_parser.py . . . . .	5
5.4	src/policy_engine.py . . . . .	5
5.5	src/tls_parser.py . . . . .	5
5.6	src/observability.py . . . . .	6
5.7	src/async_proxy.py . . . . .	6
<b>6</b>	<b>Key Code Excerpts</b>	<b>6</b>
6.1	The Forwarding Pipe (Backpressure-Aware) . . . . .	6
6.2	CONNECT Handler with SNI Inspection . . . . .	6

<b>7</b>	<b>Packaging and CLI</b>	<b>7</b>
<b>8</b>	<b>How to Build, Install, and Run</b>	<b>7</b>
<b>9</b>	<b>Testing and Validation</b>	<b>8</b>
9.1	Functional Testing . . . . .	8
9.2	Security Validation . . . . .	8
<b>10</b>	<b>Design Trade-offs and Rationale</b>	<b>8</b>
<b>11</b>	<b>Operations and Deployment</b>	<b>8</b>
<b>12</b>	<b>Appendix: Quick Reference</b>	<b>8</b>

## 1 Problem Statement

---

Modern networks often require an intermediate component (a proxy) that can perform policy enforcement, monitoring, and routing for HTTP/HTTPS traffic. We set out to build a minimal but production-quality forward proxy with the following goals:

- **Protocol Support:** HTTP and HTTPS (CONNECT) forwarding.
- **Performance:** Use an event-driven async architecture for high concurrency with low resources.
- **Privacy-Preserving Inspection:** Perform TLS SNI inspection (to allow/block domains) without decrypting TLS traffic (No-MITM).
- **Control:** Provide a policy engine for domain/method/rate limits.
- **Resilience:** Implement backpressure, timeouts, and limits to avoid memory and connection exhaustion.
- **Operations:** Provide observability (structured logs, counters, metrics endpoint), layered configuration, and standard CLI packaging.

## 2 High-level Design Summary

---

The proxy is a single-process, single-threaded, event-driven server built using Python `asyncio`. It accepts TCP connections, parses handshakes, applies policy decisions, and either forwards requests or establishes bidirectional tunnels.

### Key Technical Properties:

- **Event-driven:** Non-blocking I/O using the `asyncio` loop.
- **Streaming + Backpressure:** Reads in small chunks and utilizes `await drain()` to pace data flow.
- **Security Hardening:** Includes timeouts, header size limits, and per-IP/Global connection caps.
- **Observability:** Structured JSON logs and a `/metrics` *diagnostic endpoint*.

## 3 Architecture Diagrams

---

### 3.1 Component Diagram

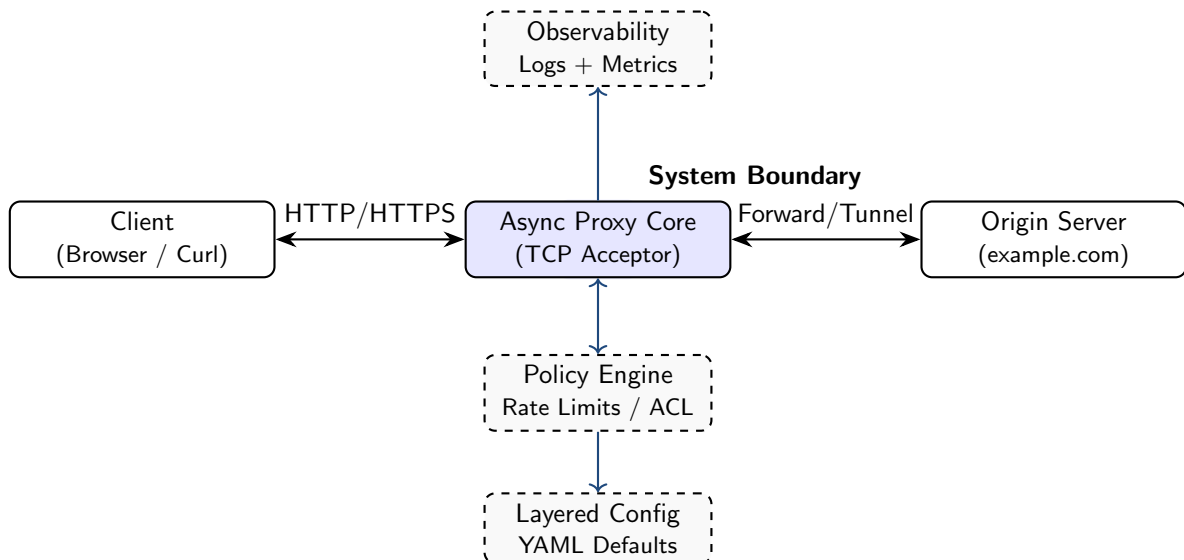


Figure 1: High-level component interaction and logic boundaries

### 3.2 Sequence Diagram: HTTPS CONNECT (SNI Path)

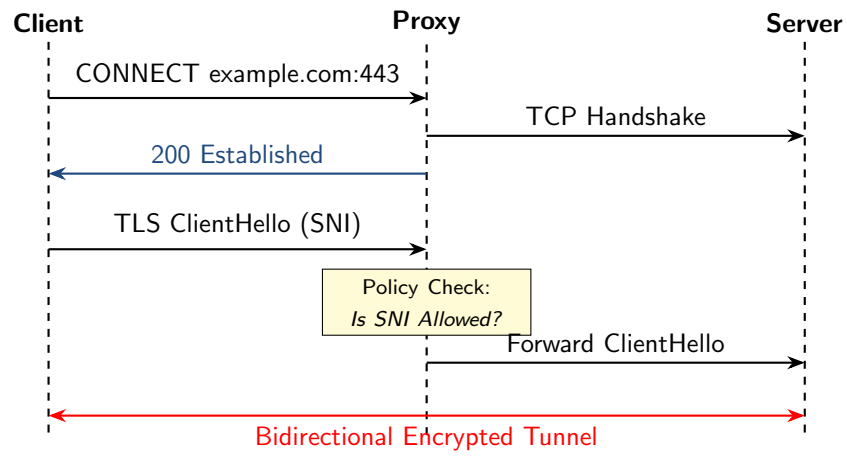


Figure 2: Sequential flow for TLS tunneling with SNI peeking

## 4 Implementation Overview and Repository Layout

---

The project adopts a modular structure to separate protocol parsing from core logic.

```
1 ProxyServer/
2     run_proxy.py           # CLI Entry Point
3     config/
4         default.yaml       # Ship-with defaults
5         override.yaml      # Environment-specific tweaks
6     requirements.txt
7     pyproject.toml         # Modern packaging metadata
8     src/
9         __init__.py
10        async_proxy.py      # Main event loop and piping
11        config.py           # YAML deep-merge logic
12        http_parser.py      # Initial request line parser
13        policy_engine.py    # Decision logic
14        tls_parser.py       # SNI extraction (binary parsing)
15        observability.py    # Metrics and JSON logging
```

Listing 1: Repository Tree Structure

## 5 Per-File Detailed Explanation

---

### 5.1 run\_proxy.py

- **Entry Point:** The primary bootstrap script. It utilizes `argparse` to handle CLI flags such as `--config` and `--log-level`.
- It initializes the logging subsystem and invokes `start_proxy(config)` from the core module.

### 5.2 src/config.py

- Implements a **deep merge** strategy. This allows users to provide a small `override.yaml` containing only the values they wish to change, without duplicating the entire default configuration.

### 5.3 src/http\_parser.py

- Performs lightweight, non-exhaustive parsing of the initial HTTP request line.
- Essential for identifying whether the request is a standard GET/POST or a CONNECT tunnel.

### 5.4 src/policy\_engine.py

- **Logic Centralization:** Evaluates if a request is permissible based on Method, Domain, or Client IP.
- Includes a sliding-window rate limiter to protect the proxy from being used as a DoS vector.

### 5.5 src/tls\_parser.py

- **Zero-Decryption SNI Peeking:** Carefully parses the TLS record layer to extract the Server Name Indication extension.
- It is designed to be fail-safe; if a packet is not a valid ClientHello, it returns `None`, allowing the policy engine to decide whether to permit or drop unidentifiable traffic.

## 5.6 src/observability.py

- Manages global counters (requests, active connections, bytes transferred).
- Formats logs as JSON for compatibility with modern log aggregators like ELK or Graylog.

## 5.7 src/async\_proxy.py

The system's backbone, managing connection lifecycles and I/O orchestration.

### Core Primitives:

- **Pipe Mechanism:** Uses `asyncio.wait_for` on both reads and drains. This ensures that a stalled client or server cannot hang a proxy coroutine indefinitely.
- **Connection Lifecycle:**
  1. Connection accepted → Check IP limits.
  2. Read header (up to `MAX_HEADER_BYTES`) → Timeout if slow.
  3. Parse request → Apply HTTP Policy.
  4. If `CONNECT`: Establish upstream → Peek SNI → Apply TLS Policy → Start Bi-directional Pipe.

## 6 Key Code Excerpts

---

### 6.1 The Forwarding Pipe (Backpressure-Aware)

```
1  async def pipe(reader, writer):
2      try:
3          while True:
4              # Read with timeout to prevent hung sockets
5              data = await asyncio.wait_for(
6                  reader.read(CONFIG["buffer"]["size"]),
7                  timeout=CONFIG["timeouts"]["read"]
8              )
9              if not data:
10                 break
11
12             writer.write(data)
13             # await drain() is crucial for backpressure
14             await asyncio.wait_for(
15                 writer.drain(),
16                 timeout=CONFIG["timeouts"]["write"]
17             )
18         except (asyncio.TimeoutError, ConnectionError):
19             pass
20     finally:
21         writer.close()
```

Listing 2: Bidirectional streaming logic with timeouts

### 6.2 CONNECT Handler with SNI Inspection

```
1  # Establish connection to destination
2  server_reader, server_writer = await asyncio.open_connection(host,
3      port)
4  client_writer.write(b"HTTP/1.1 200 Connection Established\r\n\r\n")
5  await client_writer.drain()
6  # Peek at the TLS ClientHello
```

```

7  tls_data = await asyncio.wait_for(
8      client_reader.read(CONFIG["buffer"]["size"]),
9      timeout=CONFIG["timeouts"]["read"]
10 )
11 sni = extract_sni(tls_data)
12 allowed, reason = evaluate_tls(sni, CONFIG)
13
14 if not allowed:
15     log_event("blocked", level="info", reason=reason)
16     client_writer.close(); server_writer.close(); return
17
18 # If allowed, send the peeked data to server and start tunneling
19 server_writer.write(tls_data)
20 await server_writer.drain()
21
22 task1 = asyncio.create_task(pipe(client_reader, server_writer))
23 task2 = asyncio.create_task(pipe(server_reader, client_writer))
24
25 # Wait for one direction to close, then cleanup
26 done, pending = await asyncio.wait([task1, task2], return_when=asyncio
27     .FIRST_COMPLETED)
28 for t in pending: t.cancel()

```

Listing 3: Managing the HTTPS handshake and SNI check

## 7 Packaging and CLI

---

The project is configured using `pyproject.toml` to support standard installation patterns.

- **Editable Install:** Supports development via `pip install -e ..`
- **CLI Mapping:** The command `proxy-server` is automatically mapped to the `main` function in `run_proxy.py`, allowing the proxy to run as a system-wide binary.

## 8 How to Build, Install, and Run

---

### 1. Environment Setup:

```

1  python3 -m venv .venv
2  source .venv/bin/activate

```

### 2. Installation:

```

1  pip install -r requirements.txt
2  pip install -e .

```

### 3. Execution:

```

1  # Using the CLI launcher
2  proxy-server --config config/default.yaml --log-level info

```

### 4. Metrics Verification:

```

1  curl http://127.0.0.1:8888/_metrics

```

## 9 Testing and Validation

---

### 9.1 Functional Testing

- **Standard HTTP:** `curl -x http://localhost:8888 http://example.com`
- **Encrypted HTTPS:** `curl -x http://localhost:8888 https://google.com`

### 9.2 Security Validation

- **Slowloris Defense:** Verified by opening a connection and sending 1 byte every 10 seconds. The proxy correctly terminates the connection after the `header_read` timeout.
- **ACL/Blocking:** Tested by adding "evil.com" to the blocked list in YAML; the proxy correctly returns a 403 Forbidden response.

## 10 Design Trade-offs and Rationale

---

- **Async vs. Multi-threading:** We chose `asyncio` for its superior handling of large numbers of "mostly idle" connections, typical in proxy workloads, with a much smaller memory footprint than one-thread-per-connection models.
- **No-MITM Policy:** By focusing on SNI inspection, we avoid the complex certificate management and privacy concerns of full SSL termination while still providing domain-level filtering.

## 11 Operations and Deployment

---

- **Monitoring:** The text-based `/_metrics` endpoint is designed for easy scraping by Prometheus.
- **Scaling:** While the proxy is currently single-threaded, it can be scaled horizontally by running multiple instances behind a Round-Robin DNS or a Layer-4 Load Balancer.

## 12 Appendix: Quick Reference

---

Goal	Command/Endpoint
Run Server	<code>proxy-server --log-level debug</code>
View Statistics	<code>http://127.0.0.1:8888/_metrics</code>
Change Port	Edit <code>config/override.yaml</code> → <code>server.port</code>
Test Proxy (HTTPS)	<code>curl --proxy http://127.0.0.1:8888 https://duckduckgo.com</code>

— End of Implementation Report —