

This curriculum is specifically engineered for a **Solution Architect**. In modern architecture, AWK serves as the ultimate “glue” for data pipelines, rapid prototyping of observability tools, and infrastructure-as-code (IaC) linting.

All links have been verified for December 2025.

Phase 1: The Tactical Diagnostician (Foundations)

Objective: Master the ability to extract immediate insights from raw system data and logs without installing heavy dependencies.

- **Validated Reference:** [Bruce Barnett's Grymoire AWK](#) (Classic field/record logic).
- **Architectural Concept:** Data Extraction & Pre-processing.

SA-Focused Projects:

1. **IAM Policy Audit:** Parse a tab-separated export of cloud users. Filter for those with “Admin” privileges and print their last-login date in a custom format.
 2. **Infrastructure Config Stripper:** Write a one-liner to remove comments and whitespace from HAProxy or Nginx config files for easier comparison between environments.
 3. **The CSV/JSON Field Matcher:** Given a CSV file where specific columns contain JSON strings, use field separators to extract a specific JSON key-value pair.
 4. **Resource Over-utilization Finder:** Parse a custom system report (e.g., `top` or `ps` output) and flag any process consuming more than 80% CPU using pattern matching.
 5. **Environment Variable Validator:** Write a script that checks a `.env` file for duplicate keys or missing values before a container build.
-

Phase 2: The Automation Engineer (Logic & Math)

Objective: Create scripts that perform calculations for capacity planning and resource allocation.

- **Validated Reference:** [awk.dev](#) (Official site for the **2nd Edition** of the AWK book by Aho, Kernighan, & Weinberger, released in late 2023).
- **Architectural Concept:** Reporting & Capacity Planning.

SA-Focused Projects:

1. **Cloud Bill Forecaster:** Read a CSV bill (Service, Quantity, Unit_Price). Calculate the daily spend and forecast the monthly total based on current usage.
 2. **API Latency Reporter:** Process a log of response times. Calculate the average, minimum, and maximum latency, and print a formatted summary report.
 3. **Storage Quota Estimator:** Take a list of directory sizes from `du` and convert them into GB/TB. Flag any directory exceeding a defined architectural limit.
 4. **The Transaction Tally:** Given a log of financial transactions, use logic to separate "Credit" vs "Debit" and ensure the final balance matches the header/footer records.
 5. **Multi-Log Synchronizer:** Read two logs with different timestamp formats and normalize them into a single, time-ordered view using variable math.
-

Phase 3: The Data Pipeline Architect (Associative Arrays)

Objective: Use AWK's internal memory to perform aggregations that would usually require a database.

- **Validated Reference:** [GNU AWK Manual: Arrays](#).
- **Architectural Concept:** Data Aggregation & Deduplication.

SA-Focused Projects:

1. **Unique IP/User-Agent Mapper:** Parse 10GB of Nginx logs using an associative array to count unique visitors. Output the Top 10 IP addresses by request count.
 2. **Stateful Error Tracker:** Group system errors by "Type" (e.g., Timeout, 404, 500) and keep track of the first and last time each error occurred.
 3. **The "JOIN" Simulator:** Architect a script that merges a "Server_List.csv" (ID, IP) with a "Health_Report.csv" (ID, Status) using an array to store the ID lookup.
 4. **Deduplication Engine:** Write a script that removes duplicate records in a multi-gigabyte dataset while keeping only the record with the most recent timestamp.
 5. **Cross-Region Cost Aggregator:** Given a multi-region cloud export, group costs by "Region" and "Service Type" into a multi-dimensional-like array for total architectural cost analysis.
-

Phase 4: The Observability Wizard (GAWK & Interop)

Objective: Leverage GNU AWK (GAWK) to interact with live systems, networks, and APIs.

- **Validated Reference:** [Effective AWK Programming \(Arnold Robbins\)](#) (The O'Reilly "GAWK Bible").
- **Architectural Concept:** Real-time Observability & Networking.

SA-Focused Projects:

1. **Real-Time K8s Log Watcher:** Use GAWK's two-way pipe (`|&`) to monitor `kubectl logs -f`. Trigger an alert and execute a system command if a "OOMKilled" pattern is detected.
 2. **TCP Health Checker:** Use GAWK's `/inet/tcp` to attempt connections to a list of microservices. Report which services are down without using `curl` or `telnet`.
 3. **JSON Payload Extractor:** Use GAWK's `patsplit()` or custom `RS` to parse a raw API response and extract values into environment variables for a CI/CD pipeline.
 4. **System Load Sentinel:** Use GAWK to read `/proc/loadavg` every 2 seconds. Use `strftime()` to create a timestamped CSV for Grafana ingestion.
 5. **Dynamic Traffic Router Simulation:** Write a script that reads a list of active backends and uses a round-robin logic to "assign" incoming mock requests to specific IDs.
-

Phase 5: The Enterprise Architect (Systems Design)

Objective: Build high-performance, maintainable tools that solve complex architectural challenges.

- **Validated Reference:** [GNU AWK Extensions Library](#) (For XML, CSV, and Postgres support).
- **Architectural Concept:** System Tooling & DSL (Domain Specific Languages).

SA-Focused Projects:

1. **Markdown-to-SLA Generator:** Build a tool that takes technical specifications in Markdown and generates an HTML SLA/Compliance report using AWK's string functions.
 2. **The Flat-File Metadata Store:** Architect a queryable "database" using GAWK. Support standard operations (INSERT, SELECT) to manage server metadata in a GitOps workflow.
 3. **High-Volume CSV Partitioning:** Build a script to partition a 50GB CSV file into 1GB chunks based on a specific "Tenant_ID" column, ensuring no record is split across files.
 4. **Bio-Architectural Motif Finder:** (Scientific Context) Build a tool to scan large genomic sequence files (`.fasta`) for specific patterns, a standard task in Bioinformatics solutioning.
 5. **Micro-Web Dashboard:** Use GAWK's networking to create a single-process web server that serves a "System Health" HTML page containing live CPU, RAM, and Disk metrics.
-

Summary of Solution Architect's Toolkit (2025)

Resource	URL	Purpose
Grymoire Tutorial	grymoire.com/Unix/Awk.html	Mastering field manipulation.
The New AWK Book Site	awk.dev	Modern standards (2nd Ed).
Official GAWK Manual	gnu.org/s/gawk/manual/	Deep-dive networking/internals.
Baeldung Linux Guide	baeldung.com/linux/awk-guide	Quick-start logic & formatting.
ExplainShell	explainshell.com	Visualizing complex pipelines.

Pro-Tip for Architects: In your journey, always use `gawk --profile` on your Phase 5 projects. It generates a `awkprof.out` file that shows you exactly which lines of code are execution bottlenecks, helping you design high-performance data tools.