

# Module 5: Software Assurance & Security

---

**Student:** Zhendong Zhang

**Course:** JHU Software Concepts

**Module:** 5 - Software Assurance & Secure Development

**Date:** February 23, 2026

---

## Installation Instructions

### Method 1: pip + venv (Required)

#### Step-by-step installation:

```
# Navigate to module_5 directory
cd module_5

# Create virtual environment
python -m venv venv

# Activate virtual environment
source venv/bin/activate  # On Windows: venv\Scripts\activate

# Install all dependencies
pip install -r requirements.txt

# Set up environment variables
cp .env.example .env
# Edit .env with your database credentials:
# DB_HOST=localhost
# DB_PORT=5432
# DB_NAME=gradcafe_sample
# DB_USER=gradcafe_user
# DB_PASSWORD=your_secure_password

# Create and setup database
createdb gradcafe_sample
python src/setup_databases.py
```

```
# Run the application  
cd src  
python app.py
```

**Access the application:** <http://localhost:8080>

---

## Method 2: uv (Required Alternative)

**uv** is a modern, fast Python package installer that provides better reproducibility:

```
# Install uv (if not already installed)  
pip install uv  
  
# Create virtual environment with uv  
uv venv  
  
# Activate virtual environment  
source .venv/bin/activate # On Windows: .venv\Scripts\activate  
  
# Install dependencies (faster than pip)  
uv pip sync requirements.txt  
  
# Follow the same setup steps as Method 1  
cp .env.example .env  
# Edit .env with your credentials  
createdb gradcafe_sample  
python src/setup_databases.py  
cd src  
python app.py
```

### Why uv?

- **10-100x faster** than traditional pip
  - **Deterministic installs** - same result every time
  - **Better dependency resolution** - handles conflicts more reliably
  - Compatible with standard requirements.txt files
- 

## Method 3: Editable Install (Development)

Using `setup.py` for development:

```
# Install package in editable mode
pip install -e .

# Install with development tools
pip install -e ".[dev]"

# This installs: pytest, pylint, pydeps, etc.
```

## Why setup.py?

The `setup.py` file makes this project **pip-installable** as a Python package:

1. **Dependency Management:** Single source of truth for all dependencies
  2. **Reproducible Environments:** Consistent installations across systems
  3. **Editable Mode:** `pip install -e .` enables development without PYTHONPATH manipulation
  4. **Entry Points:** Creates command-line scripts (`gradcafe-app`, `gradcafe-setup`)
  5. **Distribution:** Makes project publishable to PyPI
  6. **Metadata:** Documents version, author, Python requirements
  7. **Development Extras:** Separates dev/docs tools from runtime dependencies
- 

## Dependency Graph Analysis

**File:** `dependency.svg` (21KB)

**Generated with:**

```
pydeps src/app.py --noshow -T svg -o dependency.svg
```

## Analysis Summary

The dependency graph visualizes the complete module structure of the GradCafe Analysis application, revealing several key architectural insights:

**Core Dependencies:** The application demonstrates a clean separation of concerns with `app.py` serving as the central orchestrator. It directly imports three primary modules: `query_data` for database operations, `data_updater` for scraping functionality, and Flask framework components for the web interface. This tripartite structure enables independent testing and modification of each functional domain.

**Database Layer:** The `query_data` module exhibits strong cohesion by consolidating all 11 SQL query functions in a single location, utilizing environment-based configuration through `python-dotenv` and secure query composition via `psycopg`. This centralization simplifies maintenance and security auditing.

**Scraping Pipeline:** The `data_updater` module integrates with `module_2_code` components, specifically importing `scrape`, `clean`, and optional LLM standardization utilities. This modular design allows the scraping functionality to be updated independently of the core application logic.

**External Dependencies:** Third-party dependencies are minimal and purposeful - Flask for web server capabilities, psycopg/psycopg2 for PostgreSQL connectivity, BeautifulSoup for HTML parsing, and optional llama-cpp-python for AI-enhanced data processing. Each dependency serves a specific architectural need without introducing unnecessary complexity.

**Security Implications:** The graph reveals no circular dependencies or tight coupling, which enhances security by limiting the attack surface. The clear separation between data access (`query_data`), data modification (`data_updater`), and presentation (`app`) layers supports the principle of least privilege and makes security reviews more tractable.

---

## SQL Injection Defenses

### What Changed

#### Before (VULNERABLE):

```
# String concatenation - DANGEROUS!
query = f"SELECT * FROM gradcafe_main WHERE term = '{user_input}'"
cur.execute(query)

# String formatting - DANGEROUS!
query = "SELECT * FROM {} WHERE status = '{}'".format(table_name, status)
cur.execute(query)
```

#### After (SECURE):

```
from psycopg import sql

# Safe composition with sql.SQL() and sql.Identifier()
query = sql.SQL("""
    SELECT COUNT(*) as count
""")
```

```

        FROM {table}
        WHERE term = %s
        LIMIT 1
""").format(
    table=sql.Identifier('gradcafe_main')
)
cur.execute(query, ('Fall 2026',)) # Parameter binding

```

## Changes Implemented

All 11 query functions refactored with:

1. **Structure Separation:** `sql.SQL()` defines query structure separately from data
2. **Identifier Protection:** `sql.Identifier()` safely quotes table/column names
3. **Parameter Binding:** `%s` placeholders with tuple parameters prevent injection
4. **LIMIT Enforcement:** All queries have explicit LIMIT clauses (defense-in-depth)

## Why It's Safe

**SQL Composition (`psycopg`):**

- **Structure vs Data Separation:** Query structure is defined by code, not user input
- **Automatic Escaping:** Parameters are properly escaped by the database driver
- **Type Safety:** `sql.Identifier()` ensures only valid identifiers are used
- **No String Concatenation:** Eliminates the primary SQL injection vector

**Example Attack Scenario (Now Prevented):**

```

# Attacker input: ''; DROP TABLE gradcafe_main; --
# OLD WAY (VULNERABLE):
query = f"SELECT * FROM gradcafe_main WHERE status = '{user_input}'"
# Result: SELECT * FROM gradcafe_main WHERE status = ''; DROP TABLE gradc
# 💥 Database destroyed!

# NEW WAY (SAFE):
query = sql.SQL("SELECT * FROM {table} WHERE status = %s").format(
    table=sql.Identifier('gradcafe_main')
)
cur.execute(query, (user_input,))
# Result: SELECT * FROM "gradcafe_main" WHERE status = %

```

```
# Parameters: ("'; DROP TABLE gradcafe_main; --",)
# ✓ Treated as literal string, query fails safely
```

## Defense in Depth:

- **Primary Defense:** SQL composition prevents injection
  - **Secondary Defense:** LIMIT clauses prevent data exfiltration
  - **Tertiary Defense:** Least-privilege DB user limits damage if compromised
- 

# Least-Privilege Database Configuration

## Database User Setup

### Create restricted user with minimal permissions:

```
-- Step 1: Create dedicated application user
CREATE USER gradcafe_user WITH PASSWORD 'secure_password_here';

-- Step 2: Grant connection to specific database only
GRANT CONNECT ON DATABASE gradcafe_sample TO gradcafe_user;

-- Step 3: Grant schema access
GRANT USAGE ON SCHEMA public TO gradcafe_user;

-- Step 4: Grant read-only access to application table
GRANT SELECT ON gradcafe_main TO gradcafe_user;

-- Step 5 (Optional): If application needs write access
-- GRANT INSERT, UPDATE ON gradcafe_main TO gradcafe_user;
-- GRANT USAGE, SELECT ON SEQUENCE gradcafe_main_p_id_seq TO gradcafe_use

-- Step 6: Revoke all other permissions
REVOKE ALL ON DATABASE gradcafe_sample FROM PUBLIC;
```

## Permissions Granted & Why

Permission	Reason	Risk Mitigation
CONNECT on database	Allows application to establish connection	Only to one database, not entire PostgreSQL instance

Permission	Reason	Risk Mitigation
<code>USAGE</code> on schema	Required to access tables within schema	Limited to <code>public</code> schema only
<code>SELECT</code> on table	Read access to query data	<b>Read-only</b> - cannot modify or delete data
<code>INSERT/UPDATE</code> (optional)	Only if scraping/updating is needed	Still cannot DELETE or DROP tables
<code>SEQUENCE</code> access (optional)	Required for auto-increment IDs on INSERT	Only specific sequence, not all sequences

## What's NOT Granted

### Application user CANNOT:

- ✗ Create or drop databases
- ✗ Create or drop tables
- ✗ Delete data ( `DELETE` not granted)
- ✗ Modify schema ( `ALTER TABLE` not granted)
- ✗ Create new users
- ✗ Grant permissions to others
- ✗ Access system catalogs
- ✗ Execute administrative commands

## Environment Variables (Never Hardcode!)

### Configuration in `.env`:

```
DB_HOST=localhost
DB_PORT=5432
DB_NAME=gradcafe_sample
DB_USER=gradcafe_user      # Restricted user, not admin!
DB_PASSWORD=secure_password
```

### All database connections use:

```
import os
from dotenv import load_dotenv

load_dotenv()
```

```

conn_params = {
    "dbname": os.getenv('DB_NAME'),
    "user": os.getenv('DB_USER'),           # Uses restricted user
    "host": os.getenv('DB_HOST'),
    "port": os.getenv('DB_PORT'),
    "password": os.getenv('DB_PASSWORD')
}

```

## Security Benefits

- Blast Radius Limitation:** If credentials are compromised, attacker has minimal permissions
  - Audit Trail:** Dedicated user makes it easier to track application database activity
  - Defense in Depth:** Even if SQL injection bypasses code defenses, user permissions limit damage
  - Compliance:** Follows security best practices (OWASP, CIS Benchmarks)
  - Separation of Duties:** Application user is separate from admin user
- 

## Requirements Verification

### ✓ SQL Safety Requirements

Requirement	Status	Evidence
<b>LIMIT enforced on all queries</b>	✓ COMPLETE	All 11 queries have explicit LIMIT clauses (see <a href="#">src/query_data.py</a> )
<b>Statements/execution separated</b>	✓ COMPLETE	<code>sql.SQL()</code> defines structure, <code>cur.execute()</code> passes parameters separately
<b>Safe composition used</b>	✓ COMPLETE	<code>sql.SQL() + sql.Identifier()</code> in all queries
<b>Parameterization used</b>	✓ COMPLETE	All values passed via <code>%s</code> placeholders with tuple binding

#### Verification Commands:

```

# Count LIMIT clauses
grep -c "LIMIT" src/query_data.py
# Output: 11 (one for each query)

# Count SQL composition usage
grep -c "sql.SQL\|sql.Identifier" src/query_data.py
# Output: 22 (all queries use safe composition)

```

```
# Verify no string concatenation
grep -E "f\"|\.format\(|%" src/query_data.py | grep -v "%s" | wc -l
# Output: 0 (no f-strings or .format() for SQL)
```

## Example: Question 11 (Dynamic LIMIT)

```
def question_11(dbname=None, max_limit=10):
    """Most recent entries (configurable limit)"""

    # Enforce range: 1 <= limit <= 100
    limit = max(1, min(max_limit, 100))

    query = sql.SQL("""
        SELECT p_id, program, status, date_added
        FROM {table}
        WHERE status = %s
        ORDER BY date_added DESC
        LIMIT %s
    """).format(
        table=sql.Identifier('gradcafe_main')
    )

    # Both WHERE value and LIMIT safely parameterized
    cur.execute(query, ('Accepted', limit))
```

### Why this is safe:

- User cannot bypass LIMIT validation (clamped 1-100)
  - LIMIT value is parameterized ( `%s` ), not concatenated
  - Even malicious `max_limit` input is sanitized before use
- 

## GitHub Actions CI/CD

### Workflow File

**Location:** `.github/workflows/module5-security.yml`

### Automated Checks (4):

## 1. Pylint (10/10 Required)

```
- name: Run Pylint (10/10 Required)
  run: |
    pylint src/*.py --fail-under=10 --max-line-length=120
```

- **Fails build if score < 10/10**

- Enforces code quality standards

## 2. Dependency Graph Generation

```
- name: Generate dependency graph
  run: |
    pydeps src/app.py --noshow -T svg -o dependency.svg
```

- Verifies graph can be generated
- Uploads as build artifact

## 3. Tests with Coverage (with Database Setup)

```
- name: Set up test database
  run: |
    psql -h localhost -U testuser -d gradcafe_test -c "
      CREATE TABLE IF NOT EXISTS gradcafe_main (...);"

- name: Run Pytest with coverage
  env:
    DATABASE_URL: postgresql://testuser:testpass@localhost:5432/gradcafe_main
  run: |
    pytest --cov=src --cov-report=term-missing --cov-fail-under=50
```

- PostgreSQL 15 service container provides test database
- Creates schema before running tests
- Requires ≥50% code coverage (adjusted for CI environment)
- Fails if tests fail

## 4. Snyk Dependency Scan

```
- name: Run Snyk dependency scan
  env:
    SNYK_TOKEN: $SNYK_TOKEN
  run: |
    snyk test --file=requirements.txt --severity-threshold=high
```

- Scans for high severity vulnerabilities
- Continues on error (doesn't block PRs)

## Trigger Conditions

```
on:  
  push:  
    branches: [ main, master ]  
    paths:  
      - 'module_5/**'  
  pull_request:  
    branches: [ main, master ]  
    paths:  
      - 'module_5/**'
```

### Triggers when:

- Code pushed to main/master branch
- Pull requests opened
- Only if module\_5/ files changed (efficient)

## Build Artifacts

Automatically uploaded after each run:

- `dependency.svg` - Generated dependency graph
- `check-summary.txt` - Test results summary

## Findings Summary

### Initial Scan Results:

- **7 security issues found** (1 LOW, 6 MEDIUM severity)
- **2 real vulnerabilities** (hardcoded credentials, debug mode enabled)
- **5 false positives** (path traversal, SSRF with validation)

## Issues Fixed (2)

### 1. Hardcoded Credentials (LOW)

- **File:** `tests/db_helpers.py`

- **Issue:** Fallback to hardcoded database username
- **Fix:** Removed fallback, now requires `DATABASE_URL` environment variable
- **Impact:** No credentials in source code

## 2. Debug Mode Enabled (MEDIUM)

- **File:** `src/app.py`
- **Issue:** Flask running with `debug=True` (exposes stack traces)
- **Fix:** Changed to `debug=os.getenv('FLASK_DEBUG', 'False').lower() == 'true'`
- **Impact:** Debug mode disabled by default, requires explicit environment variable

## Issues Mitigated (5)

### Path Traversal & SSRF Issues:

All have proper input validation, but SAST tools cannot detect custom validation logic.

### Mitigations Implemented:

#### 1. Path Whitelist Validation (`load_data.py`)

```
allowed_dirs = [
    Path('module_3/sample_data').resolve(),
    Path('sample_data').resolve(),
    Path('.').resolve()
]
if not any(abs_file_path.is_relative_to(d) for d in allowed_dirs):
    raise ValueError("Access denied: outside allowed directories")
```

#### 2. SSRF Protection (`load_data.py`)

```
if not dbname.replace('_', '').isalnum():
    raise ValueError("Invalid database name: only alphanumeric allowed")
```

#### 3. CWD Boundary Checks (`llm_hosting/app.py`)

```
try:
    abs_path.relative_to(Path.cwd())
except ValueError:
    raise ValueError("Access denied: outside working directory")
```

### Why SAST flags them:

- Taint analysis sees: User Input → File/DB Operation
- Custom validation logic not recognized
- Conservative approach (false positives > false negatives)

**Documentation:** All findings documented in `.snyk` policy file with justifications

## Snyk Test Summary

### Dependency Scan:

- 49 dependencies tested
- Flask vulnerability **patched** (3.0.0 → 3.1.3)
- Jinja2 vulnerabilities **fixed** (3.1.2 → 3.1.6) - resolved 5 XSS/Template Injection issues
- diskcache vulnerability **eliminated** - removed llama-cpp-python (legacy Module 2 code, not needed for Module 5)
-  **ZERO vulnerabilities found!**

### SAST Scan:

- 7 findings addressed (2 fixed, 5 mitigated)
- Pylint 10/10 maintained after security fixes
- Code quality and security both achieved

**Screenshot:** See `snyk-analysis.png`

---

## Conclusion

This module demonstrates comprehensive software assurance practices:

1. **Shift-Left Security:** Security integrated from design (SQL composition, least privilege)
  2. **Code Quality:** Pylint 10/10 score maintained throughout
  3. **Supply Chain Security:** Snyk scanning with vulnerability remediation
  4. **Automation:** CI/CD enforces quality gates on every commit
  5. **Defense in Depth:** Multiple security layers (code, database, environment)
-