

# /sentry-security-review

**Source:** `~/ .claude/skills/sentry-security-review/SKILL.md`

---

name: security-review description:  
Security code review for vulnerabilities.  
Use when asked to "security review",  
"find vulnerabilities", "check for security  
issues", "audit security", "OWASP  
review", or review code for injection,  
XSS, authentication, authorization,  
cryptography issues. Provides  
systematic review with confidence-  
based reporting. allowed-tools: Read  
Grep Glob Bash Task license:  
LICENSE

## Security Review Skill

Identify exploitable security vulnerabilities in code. Report only **HIGH CONFIDENCE** findings—clear vulnerable patterns with attacker-controlled input.

# Scope: Research vs. Reporting

## CRITICAL DISTINCTION:

- **Report on:** Only the specific file, diff, or code provided by the user
- **Research:** The ENTIRE codebase to build confidence before reporting

Before flagging any issue, you **MUST** research the codebase to understand:

- Where does this input actually come from? (Trace data flow)
- Is there validation/sanitization elsewhere?
- How is this configured? (Check settings, config files, middleware)
- What framework protections exist?

**Do NOT report issues based solely on pattern matching.** Investigate first, then report only what you're confident is exploitable.

# Confidence Levels

Level	Criteria	Action
<b>HIGH</b>	Vulnerable pattern + attacker-controlled input confirmed	<b>Report</b> with severity
<b>MEDIUM</b>	Vulnerable pattern, input source unclear	<b>Note</b> as "Needs verification"
<b>LOW</b>	Theoretical, best practice, defense-in-depth	<b>Do not report</b>

# Do Not Flag

# General Rules

- Test files (unless explicitly reviewing test security)
- Dead code, commented code, documentation strings
- Patterns using **constants** or **server-controlled configuration**
- Code paths that require prior authentication to reach (note the auth requirement instead)

# Server-Controlled Values (NOT Attacker-Controlled)

These are configured by operators, not controlled by attackers:

Source	Example	Why It's Safe
Django settings	<code>settings.API_URL</code> , <code>settings.ALLOWED_HOSTS</code>	Set via config/env at deployment
Environment variables	<code>os.environ.get('DATABASE_URL')</code>	Deployment configuration
Config files	<code>config.yaml</code> , <code>app.config['KEY']</code>	Server-side files
Framework constants	<code>django.conf.settings.*</code>	Not user-modifiable
Hardcoded values	<code>BASE_URL = "https://api.internal"</code>	Compile-time constants

## SSRF Example - NOT a vulnerability:

```
# SAFE: URL comes from Django settings (server-controlled)
response = requests.get(f"{settings.SSEER_AUTOFIX_URL}{path}")
```

## SSRF Example - IS a vulnerability:

```
# VULNERABLE: URL comes from request (attacker-controlled)
response = requests.get(request.GET.get('url'))
```

# Framework-Mitigated Patterns

Check language guides before flagging. Common false positives:

Pattern	Why It's Usually Safe
Django <code>{{ variable }}</code>	Auto-escaped by default
React <code>{variable}</code>	Auto-escaped by default
Vue <code>{{ variable }}</code>	Auto-escaped by default
<code>User.objects.filter(id=input)</code>	ORM parameterizes queries
<code>cursor.execute("...%S", (input,))</code>	Parameterized query
<code>innerHTML = "&lt;b&gt;Loading...&lt;/b&gt;"</code>	Constant string, no user input

## Only flag these when:

- Django: `{{ var|safe }}`, `{% autoescape off %}`, `mark_safe(user_input)`

- React: `dangerouslySetInnerHTML={{__html: userInput}}`
- Vue: `v-html="userInput"`
- ORM: `.raw()`, `.extra()`, `RawSQL()` with string interpolation

# Review Process

## 1. Detect Context

What type of code am I reviewing?

Code Type	Load These References
API endpoints, routes	<code>authorization.md</code> , <code>authentication.md</code> , <code>injection.md</code>
Frontend, templates	<code>xss.md</code> , <code>csrf.md</code>
File handling, uploads	<code>file-security.md</code>
Crypto, secrets, tokens	<code>cryptography.md</code> , <code>data-protection.md</code>
Data serialization	<code>deserialization.md</code>
External requests	<code>ssrf.md</code>
Business workflows	<code>business-logic.md</code>
GraphQL, REST design	<code>api-security.md</code>
Config, headers, CORS	<code>misconfiguration.md</code>
CI/CD, dependencies	<code>supply-chain.md</code>
Error handling	<code>error-handling.md</code>
Audit, logging	<code>logging.md</code>

## 2. Load Language Guide

Based on file extension or imports:

Indicators	Guide
<code>.py</code> , <code>django</code> , <code>flask</code> , <code>fastapi</code>	<code>languages/python.md</code>
<code>.js</code> , <code>.ts</code> , <code>express</code> , <code>react</code> , <code>vue</code> , <code>next</code>	<code>languages/javascript.md</code>
<code>.go</code> , <code>go.mod</code>	<code>languages/go.md</code>
<code>.rs</code> , <code>Cargo.toml</code>	<code>languages/rust.md</code>
<code>.java</code> , <code>spring</code> , <code>@Controller</code>	<code>languages/java.md</code>

## 3. Load Infrastructure Guide (if applicable)

File Type	Guide
Dockerfile, .dockerignore	infrastructure/docker.md
K8s manifests, Helm charts	infrastructure/kubernetes.md
.tf, Terraform	infrastructure/terraform.md
GitHub Actions, .gitlab-ci.yml	infrastructure/ci-cd.md
AWS/GCP/Azure configs, IAM	infrastructure/cloud.md

## 4. Research Before Flagging

For each potential issue, research the codebase to build confidence:

- Where does this value actually come from? Trace the data flow.
- Is it configured at deployment (settings, env vars) or from user input?
- Is there validation, sanitization, or allowlisting elsewhere?
- What framework protections apply?

Only report issues where you have HIGH confidence after understanding the broader context.

## 5. Verify Exploitability

For each potential finding, confirm:

Is the input attacker-controlled?

Attacker-Controlled (Investigate)	Server-Controlled (Usually Safe)
request.GET, request.POST, request.args	settings.X, app.config['X']
request.json, request.data, request.body	os.environ.get('X')
request.headers (most headers)	Hardcoded constants
request.cookies (unsigned)	Internal service URLs from config
URL path segments: /users/<id>/	Database content from admin/system
File uploads (content and names)	Signed session data
Database content from other users	Framework settings
WebSocket messages	

Does the framework mitigate this?

- Check language guide for auto-escaping, parameterization

- Check for middleware/decorators that sanitize

### Is there validation upstream?

- Input validation before this code
- Sanitization libraries (DOMPurify, bleach, etc.)

## 6. Report HIGH Confidence Only

Skip theoretical issues. Report only what you've confirmed is exploitable after research.

## Severity Classification

Severity	Impact	Examples
<b>Critical</b>	Direct exploit, severe impact, no auth required	RCE, SQL injection to data, auth bypass, hardcoded secrets
<b>High</b>	Exploitable with conditions, significant impact	Stored XSS, SSRF to metadata, IDOR to sensitive data
<b>Medium</b>	Specific conditions required, moderate impact	Reflected XSS, CSRF on state-changing actions, path traversal
<b>Low</b>	Defense-in-depth, minimal direct impact	Missing headers, verbose errors, weak algorithms in non-critical context

## Quick Patterns Reference

### Always Flag (Critical)

```
eval(user_input)           # Any language
exec(user_input)          # Any language
pickle.loads(user_data)   # Python
yaml.load(user_data)      # Python (not safe_load)
unserialize($user_data)   # PHP
deserialize(user_data)    # Java ObjectInputStream
shell=True + user_input   # Python subprocess
child_process.exec(user)  # Node.js
```

# Always Flag (High)

```
innerHTML = userInput           # DOM XSS
dangerouslySetInnerHTML={user} # React XSS
v-html="userInput"             # Vue XSS
f"SELECT * FROM x WHERE {user}" # SQL injection
`SELECT * FROM x WHERE ${user}` # SQL injection
os.system(f"cmd {user_input}")  # Command injection
```

# Always Flag (Secrets)

```
password = "hardcoded"
api_key = "sk-..."
AWS_SECRET_ACCESS_KEY = "..."
private_key = "-----BEGIN"
```

# Check Context First (MUST Investigate Before Flagging)

```
# SSRF - ONLY if URL is from user input, NOT from settings/config
requests.get(request.GET['url']) # FLAG: User-controlled URL
requests.get(settings.API_URL)  # SAFE: Server-controlled config
requests.get(f"{settings.BASE}/{x}") # CHECK: Is 'x' user input?

# Path traversal - ONLY if path is from user input
open(request.GET['file'])       # FLAG: User-controlled path
open(settings.LOG_PATH)        # SAFE: Server-controlled config
open(f"{BASE_DIR}/{filename}") # CHECK: Is 'filename' user input?

# Open redirect - ONLY if URL is from user input
redirect(request.GET['next'])   # FLAG: User-controlled redirect
redirect(settings.LOGIN_URL)   # SAFE: Server-controlled config

# Weak crypto - ONLY if used for security purposes
hashlib.md5(file_content)     # SAFE: File checksums, caching
hashlib.md5(password)        # FLAG: Password hashing
random.random()               # SAFE: Non-security uses (UI, sampling)
random.random() for token     # FLAG: Security tokens need secrets module
```

# Output Format

```
## Security Review: [File/Component Name]

### Summary
- **Findings**: X (Y Critical, Z High, ...)
- **Risk Level**: Critical/High/Medium/Low
- **Confidence**: High/Mixed

### Findings

#### [VULN-001] [Vulnerability Type] (Severity)
- **Location**: `file.py:123`
- **Confidence**: High
- **Issue**: [What the vulnerability is]
- **Impact**: [What an attacker could do]
- **Evidence**:
  ```python
  [Vulnerable code snippet]
```

- **Fix:** [How to remediate]

## Needs Verification

### [VERIFY-001] [Potential Issue]

- **Location:** `file.py:456`
- **Question:** [What needs to be verified]

```
If no vulnerabilities found, state: "No high-confidence vulnerabilities identified."
```

```
---
```

```
## Reference Files
```

```
### Core Vulnerabilities (`references/`)
```

```
| File | Covers |
```

```
|-----|-----|
```

```
| `injection.md` | SQL, NoSQL, OS command, LDAP, template injection |
| `xss.md` | Reflected, stored, DOM-based XSS |
| `authorization.md` | Authorization, IDOR, privilege escalation |
| `authentication.md` | Sessions, credentials, password storage |
| `cryptography.md` | Algorithms, key management, randomness |
| `deserialization.md` | Pickle, YAML, Java, PHP deserialization |
| `file-security.md` | Path traversal, uploads, XXE |
| `ssrf.md` | Server-side request forgery |
| `csrf.md` | Cross-site request forgery |
| `data-protection.md` | Secrets exposure, PII, logging |
| `api-security.md` | REST, GraphQL, mass assignment |
| `business-logic.md` | Race conditions, workflow bypass |
| `modern-threats.md` | Prototype pollution, LLM injection, WebSocket |
| `misconfiguration.md` | Headers, CORS, debug mode, defaults |
| `error-handling.md` | Fail-open, information disclosure |
| `supply-chain.md` | Dependencies, build security |
| `logging.md` | Audit failures, log injection |
```

### ### Language Guides (`languages/`)

- `python.md` - Django, Flask, FastAPI patterns
- `javascript.md` - Node, Express, React, Vue, Next.js
- `go.md` - Go-specific security patterns
- `rust.md` - Rust unsafe blocks, FFI security
- `java.md` - Spring, Java EE patterns

### ### Infrastructure (`infrastructure/`)

- `docker.md` - Container security
- `kubernetes.md` - K8s RBAC, secrets, policies
- `terraform.md` - IaC security
- `ci-cd.md` - Pipeline security
- `cloud.md` - AWS/GCP/Azure security

---

Revision #5

Created 2026-02-18 08:40:01 UTC by John

Updated 2026-06-21 20:01:07 UTC by John