# Analyse statique de code avec Semgrep

**SSTIC 2023**

Claudio Merloni (claudio@semgrep.com)
https://bit.ly/semgrep_sstic2023

Semgrep

# Agenda

1. Introduction to Semgrep OSS
   a. Philosophy
   b. Architecture and features
   c. Getting started
   d. Alternatives
2. Finding vulnerabilities
3. Enforcing secure defaults
4. Exploring a code base
5. References

# About me

**Claudio Merloni**
Security Research Manager @ Semgrep
claudio@semgrep.com
claudiomerloni

Proudly stuck in static analysis since 2008
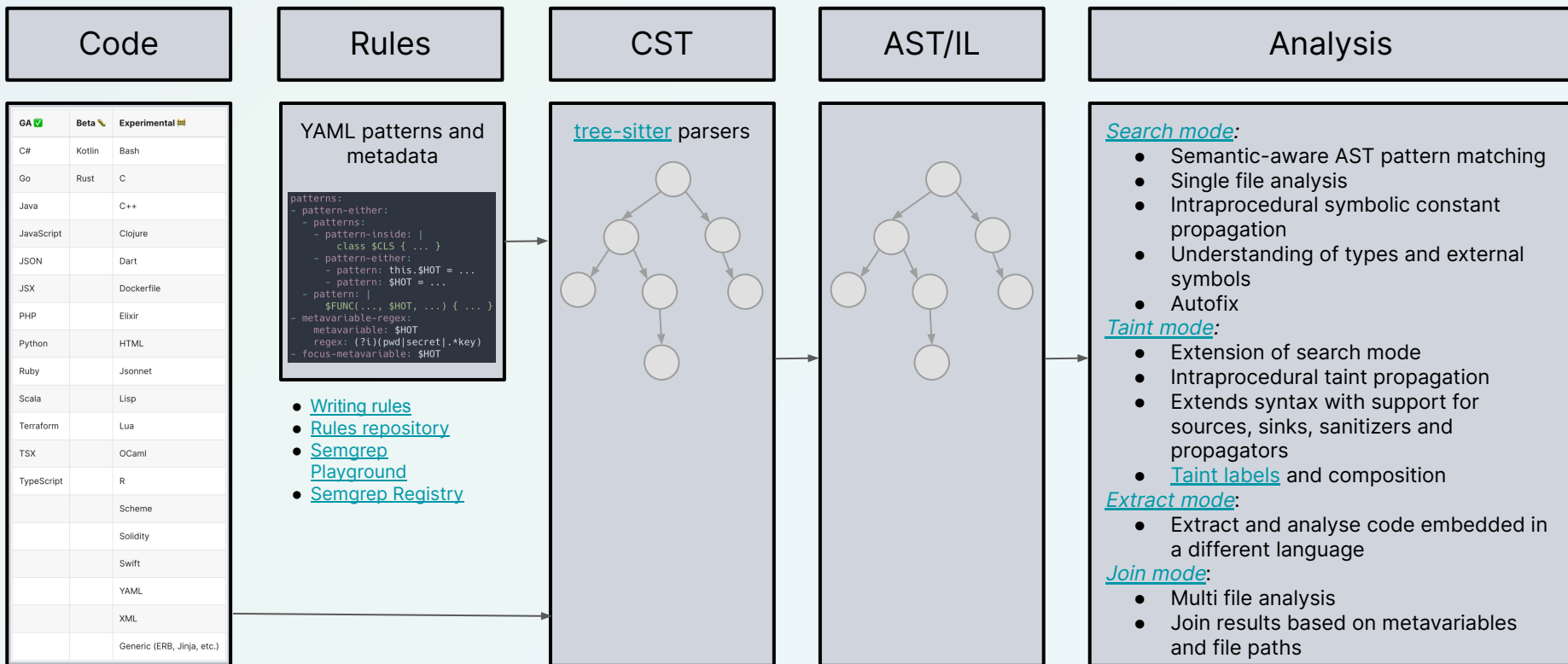Formerly: Synopsys, Fortify

# Semgrep OSS

# Philosophy

*TL;DR*

- Free & OSS
- Local & self-contained
- Runs everywhere, analyses everything
- User-friendly: get started using it and writing rules in minutes

https://semgrep.dev/docs/contributing/semgrep-philosophy/

# Analysis architecture

| Code | Rules | CST | AST/IL | Analysis |
|------|-------|-----|--------|----------|

**Code**

| GA ✅ | Beta 🔧 | Experimental 🚧 |
|-------|---------|------------------|
| C# | Kotlin | Bash |
| Go | Rust | C |
| Java | | C++ |
| JavaScript | | Clojure |
| JSON | | Dart |
| JSX | | Dockerfile |
| PHP | | Elixir |
| Python | | HTML |
| Ruby | | Jsonnet |
| Scala | | Lisp |
| Terraform | | Lua |
| TSX | | OCaml |
| TypeScript | | R |
| | | Scheme |
| | | Solidity |
| | | Swift |
| | | YAML |
| | | XML |
| | | Generic (ERB, Jinja, etc.) |

Supported Languages

**Rules**

YAML patterns and metadata

```
patterns:
- pattern-either:
  - patterns:
    - pattern-inside: |
        class $CLS { ... }
    - pattern-either:
      - pattern: this.$HOT = ...
      - pattern: $HOT = ...
  - pattern: |
      $FUNC(..., $HOT, ...) { ... }
- metavariable-regex:
    metavariable: $HOT
    regex: (?i)(pwd|secret|.*key)
- focus-metavariable: $HOT
```

- Writing rules
- Rules repository
- Semgrep Playground
- Semgrep Registry

**CST**

tree-sitter parsers

**AST/IL**

**Analysis**

*Search mode*:
- Semantic-aware AST pattern matching
- Single file analysis
- Intraprocedural symbolic constant propagation
- Understanding of types and external symbols
- Autofix

*Taint mode*:
- Extension of search mode
- Intraprocedural taint propagation
- Extends syntax with support for sources, sinks, sanitizers and propagators
- Taint labels and composition

*Extract mode*:
- Extract and analyse code embedded in a different language

*Join mode*:
- Multi file analysis
- Join results based on metavariables and file paths

# Getting started

- *Prerequisites*: Python, WSL on Windows
- *Installation*
    - brew install semgrep
    - pip install semgrep (or pipx install semgrep)
- *Usage*
    - semgrep --config=<rules> <my target folder>
- *Alternatively*
    - docker run --rm -v "${PWD}:/src" returntocorp/semgrep semgrep --config=<rules>


- Can be used in CLI, pre-commit, CI workflows, etc.
- See https://semgrep.dev/docs/getting-started/

# Alternatives

- *CodeQL*
  - Engine is not OSS
  - Open to [contributions](#) to queries and libraries
  - Free to use only on public GitHub repositories
  - [Supports](#) C, C++, C#, Go, Java, Kotlin, JS/TS, Python, Ruby
  - Requires full "buildable" code base
  - Query syntax has rather high learning curve
- *Weggli*
  - Written in Rust, doesn't depend on something like a Python installation
  - Parsing based on tree-sitter
  - Focused on C and C++
  - Doesn't require buildable code base
  - Patterns are very close to C/C++ code
  - [Stay here to hear more about it from Kevin Denis!](#)
- ... and many others! *Gosec, Brakeman, Bandit, Gitleaks, ...*

# Finding vulnerabilities

# Audit for SSRF

```java
protected AttackResult furBall(String url) {
    if (url.matches("http://ifconfig\\.pro")) {
        String html;
        try (InputStream in = new URL(url)
                .openStream()) { /* ... */ }
```

```yaml
patterns:
- pattern-either:
  - patterns:
    - pattern-inside: |
        $RETURN $METHOD(..., String $ARG, ...) {
            ...
        }
    - pattern: new java.net.URL($ARG)
  - patterns:
    - pattern: new java.net.URL(...)
    - pattern-not: new java.net.URL("...")
```

Try it out: https://semgrep.dev/s/J4jw
Similar vulnerability, different audit technique: https://blog.doyensec.com/2023/03/16/ssrf-remediation-bypass.html

# Hunting for SQLi

```typescript
const someSQLQuery = (req: express.Request): string | null => {
  const pool=new Pool(a)
  pool.query("INSERT INTO foo (p, desc, s) VALUES ('"+
              req.body.p + "', '" +
              req.body.d + "', 'Ok');");
};
```

```yaml
mode: taint
pattern-sources:
  - pattern: |
      ($REQ: express.Request).body
pattern-sanitizers:
  - pattern: parseInt(...)
pattern-sinks:
  - patterns:
      - pattern: $POOL.query($VALUE,...)
      - focus-metavariable: $VALUE
```

Try it out: https://semgrep.dev/playground/s/2oyj

# Hunting for hardcoded secrets

```yaml
1 mode: taint
2 pattern-sources:
3   - pattern: |
4       {..., password:'...', ...}
5 pattern-sinks:
6   - patterns:
7       - pattern-inside: |
8           $DB = require('$LIB')
9           ...
10      - metavariable-regex:
11          metavariable: $LIB
12          regex: (somedb|otherdb)
13      - pattern: |
14          $DB.connect($FOO)
15      - focus-metavariable: $FOO
```

```javascript
1  let fuu = "password"
2  let bar = {
3  host: "test.test.ap-east-1.rds.amazonaws.com",
4  user: "newuser",
5  password: fuu,
6  database: "test",
7  port: "3306",
8  }
9  // ruleid: more-db-hardcoded-secret
10 var conn = db.connect(bar);
```

Try it out: https://semgrep.dev/s/vypW

# Enforcing secure defaults

# Stop leaking sensitive data

```python
1  from django.db import Model, Column, String
2
3  class Token(Model):
4      bad_token = Column(String, nullable=False, unique=True, index=True)
```

```yaml
1  patterns:
2    - pattern: |
3        $COLUMN = django.db.Column(django.db.String, ...)
4    - metavariable-regex:
5        metavariable: $COLUMN
6        regex: (?i)(.*(?<![A-Za-z])(token|key|email|secret|password)(?![A-RT-Za-rt-z]).*)
```

Try it out: https://semgrep.dev/s/0klB

# And code conventions in general!

- Validate format of structured data files (JSON, YAML, etc.)
  - We use it for our own Semgrep rules!
- Enforcing naming conventions (classes, API endpoints, etc.)
- Usage of forbidden/unsafe libraries and APIs

Example: https://www.fabianzeindl.com/posts/business-information-server

# Exploring a code base

# Discovering web application routes

```java
1  @Controller
2  @RequestMapping("/api/test")
3  public class ExampleController {
4
5      @RequestMapping(method = RequestMethod.GET)
6      @Authorize(Permissions.ADMIN)
7      @ResponseBody
8      public ResponseEntity<Map<String, Object>> list() {
9          return new ResponseEntity<>(result, HttpStatus.OK);
10     }
11
12     @RequestMapping(method = RequestMethod.GET)
13     @ResponseBody
14     public ResponseEntity<Map<String, Object>> unauth() {
15         return new ResponseEntity<>(result, HttpStatus.OK);
16     }
17
18     @RequestMapping(value = "/{name}", method = RequestMethod.POST)
19     @Authorize(Permissions.USER)
20     @ResponseBody
21     public ResponseEntity<Map<String, Object>> load(@PathVariable final String name) throws APIException {
22         return new ResponseEntity<>(result, HttpStatus.OK);
23     }
24 }
```

# Discovering web application routes

```
1  patterns:
2    - pattern-inside: |
3        @Controller
4        public class $CONTROLLER { ... }
5    - pattern-either:
6        - pattern: |
7            @RequestMapping(method = $HTTPMETHOD)
8            @Authorize($AUTHZ)
9            public $RETURNTYPE $METHOD(...) { ...}
10       - pattern: |
11           @RequestMapping(method = $HTTPMETHOD, value = $PATH)
12           @Authorize($AUTHZ)
13           public $RETURNTYPE $METHOD(...) { ...}
14       - pattern: |
15           @RequestMapping(method = $HTTPMETHOD)
16           public $RETURNTYPE $METHOD(...) { ...}
```

Try it out: https://semgrep.dev/s/WdqL

# References

# References

- Semgrep source code: https://github.com/returntocorp/semgrep
- Registry: https://semgrep.dev/explore
  - Community rules repository: https://github.com/returntocorp/semgrep-rules
- Rule writing tutorial: https://semgrep.dev/learn
- Docs: https://semgrep.dev/docs/
- Community Slack: https://r2c.dev/slack

- https://blog.trailofbits.com/2021/11/08/discovering-goroutine-leaks-with-semgrep/
- https://blog.includesecurity.com/2021/07/customizing-semgrep-rules-for-flask-django/
- https://notsosecure.com/semgrep-practical-introduction
- https://blog.aquia.io/blog/2022-02-18-semgrep-cdk/
- https://blog.doyensec.com/2023/03/16/ssrf-remediation-bypass.html
- https://www.fabianzeindl.com/posts/business-information-server

# Merci!

# More examples

# SQL Injection with *taint labels* - The code

```
1  package main
2
3  import (
4      "fmt"
5      "net/http"
6
7      "github.com/jackc/pgx/v5"
8  )
9
10 func bad(w http.ResponseWriter, req *http.Request) {
11     pgxConfig := pgx.ConnConfig{
12         Host:     "localhost",
13         Database: "quetest",
14         User:     "quetest",
15     }
16     pgxConnPoolConfig := pgx.ConnPoolConfig{pgxConfig, 3, nil}
17     conn, err := pgx.NewConnPool(pgxConnPoolConfig)
18     if err != nil {
19         log.Fatal(err)
20     }
21     query = "SELECT name FROM users WHERE age=" + req.FormValue("age")
22     // ruleid: pgx-sqli-example
23     rows, err := conn.Query(context.Background(), query)
24 }
```

Focus on a specific library

We know where user data comes from

The usual suspect …

Try it out: https://semgrep.dev/s/klD1

# SQL Injection with *taint labels* - Database library

```
1  - patterns:
2     - pattern-inside: |
3         import "$IMPORT"
4         ...
5     - metavariable-regex:
6         metavariable: $IMPORT
7         regex: (.*jackc\/pgx\/v(4|5).*)
8  label: IMPORTPGX
```

```
1  import (
2      "fmt"
3      "net/http"
4
5      "github.com/jackc/pgx/v5"
6  )
```

# SQL Injection with *taint labels* - User input

```yaml
1  - patterns:
2      - pattern-inside: |
3          import "net/http"
4          ...
5      - pattern-either:
6          - pattern: |
7              ($REQ : http.Request).$FIELD
8          - pattern: |
9              ($REQ : *http.Request).$FIELD
10     - metavariable-regex:
11         metavariable: $FIELD
12         regex: ^(FormValue)$
13 label: USERINPUT
```

```go
1  import (
2      "fmt"
3      "net/http"
4
5      "github.com/jackc/pgx/v5"
6  )
7
8  func bad(w http.ResponseWriter, req *http.Request) {
9      pgxConfig := pgx.ConnConfig{
10         Host:     "localhost",
11         Database: "quetest",
12         User:     "quetest",
13     }
14     pgxConnPoolConfig := pgx.ConnPoolConfig{pgxConfig, 3, nil}
15     conn, err := pgx.NewConnPool(pgxConnPoolConfig)
16     if err != nil {
17         log.Fatal(err)
18     }
19     query = "SELECT name FROM users WHERE age=" + req.FormValue("age")
```

# SQL Injection with *taint labels* - Putting it together

```
1  pattern-sinks:
2    - patterns:
3        - patterns:
4            - pattern-either:
5                - pattern-inside: |
6                    $DB := pgx.$CONNECT(...)
7                    ...
8                - pattern-inside: |
9                    $DB, ... := pgx.$CONNECT(...)
10                   ...
11           - metavariable-regex:
12               metavariable: $CONNECT
13               regex: ^(NewConnPool)$
14         - pattern: $DB. ... .$METHOD($CTX, $QUERY, ...)
15         - metavariable-regex:
16             metavariable: $METHOD
17             regex: ^(Query|QueryRow)$
18         - focus-metavariable: $QUERY
19     requires: IMPORTPGX and USERINPUT and CONCAT
```

Focus on a block of code and capture the database connection

Call on a method of the captured variable

The query is our sink ...

... only if source constraints are met!

# Forgotten debugging code

```
1 pattern-either:
2   - pattern: pdb.$X(...)
3   - pattern: pdb.Pdb.$X(...)
```

```
1 import pdb as db
2
3 def foo():
4     # ruleid:pdb-remove
5     db.set_trace()
6     # ok:pdb-remove
7     a = "apple"
8     #ok:pdb-remove
9     db = "the string, not the library"
10    #ok:pdb-remove
11    pdb = "also a string"
12    # ruleid:pdb-remove
13    pdb.Pdb.set_trace()
14    # ruleid:pdb-remove
15    db.Pdb.set_trace(...)
```

Try it out: https://semgrep.dev/r/python.lang.correctness.pdb.pdb-remove