



jÄk: Using Dynamic Analysis to Crawl and Test Modern Web Applications

Giancarlo Pellegrino⁽¹⁾, Constantin Tschürtz⁽²⁾, Eric Bodden⁽²⁾, and Christian Rossow⁽¹⁾

18th International Symposium on Research in Attacks, Intrusions and Defenses

November 3rd, Kyoto, Japan

⁽¹⁾ CISPA, Saarland University, Germany

⁽²⁾ Fraunhofer SIT / TU Darmstadt, Germany

Web Application Scanners



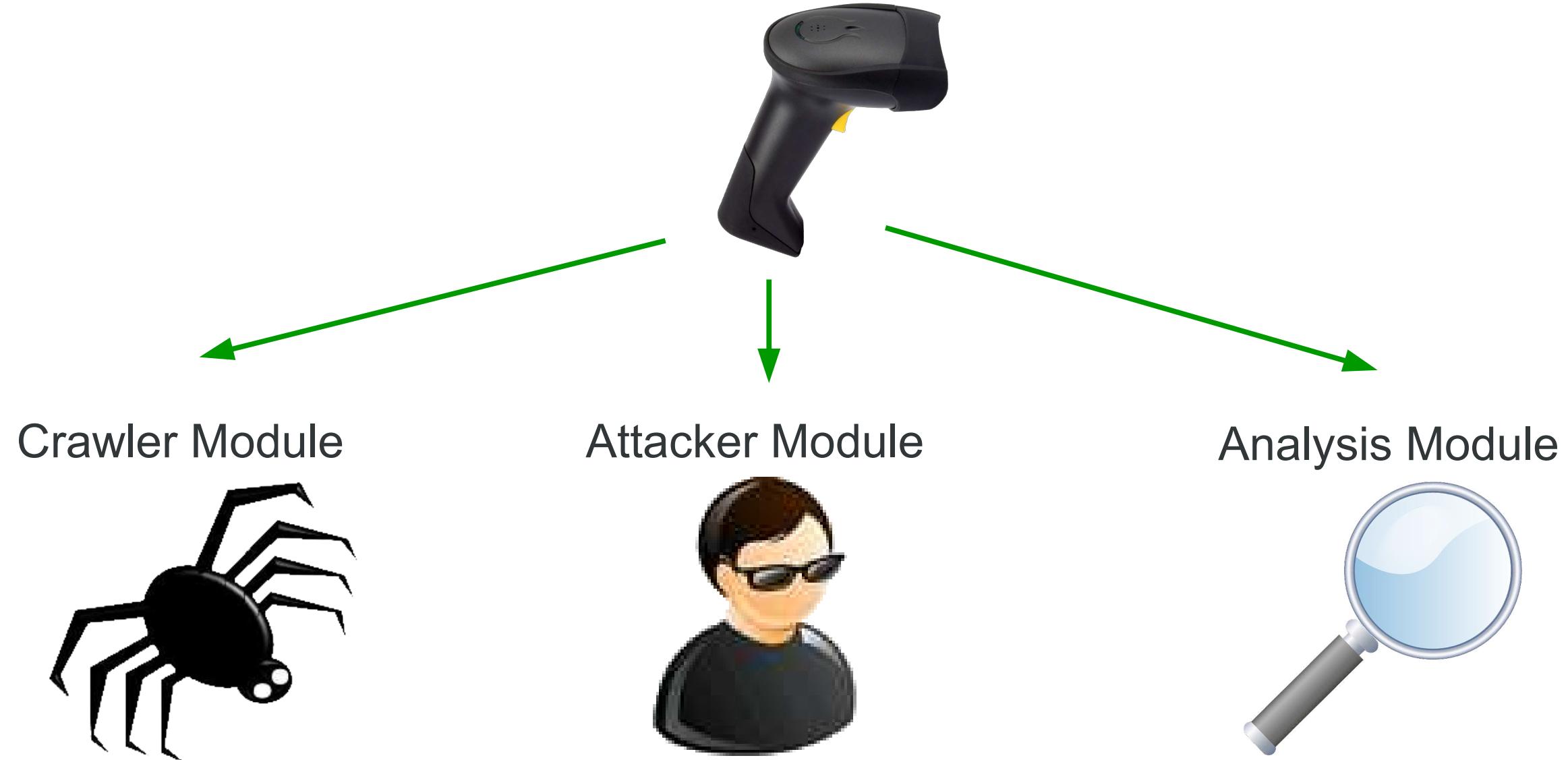
- (Semi-)automated security testing tools
- Follow a dynamic and black-box testing approach

Web Application Scanners

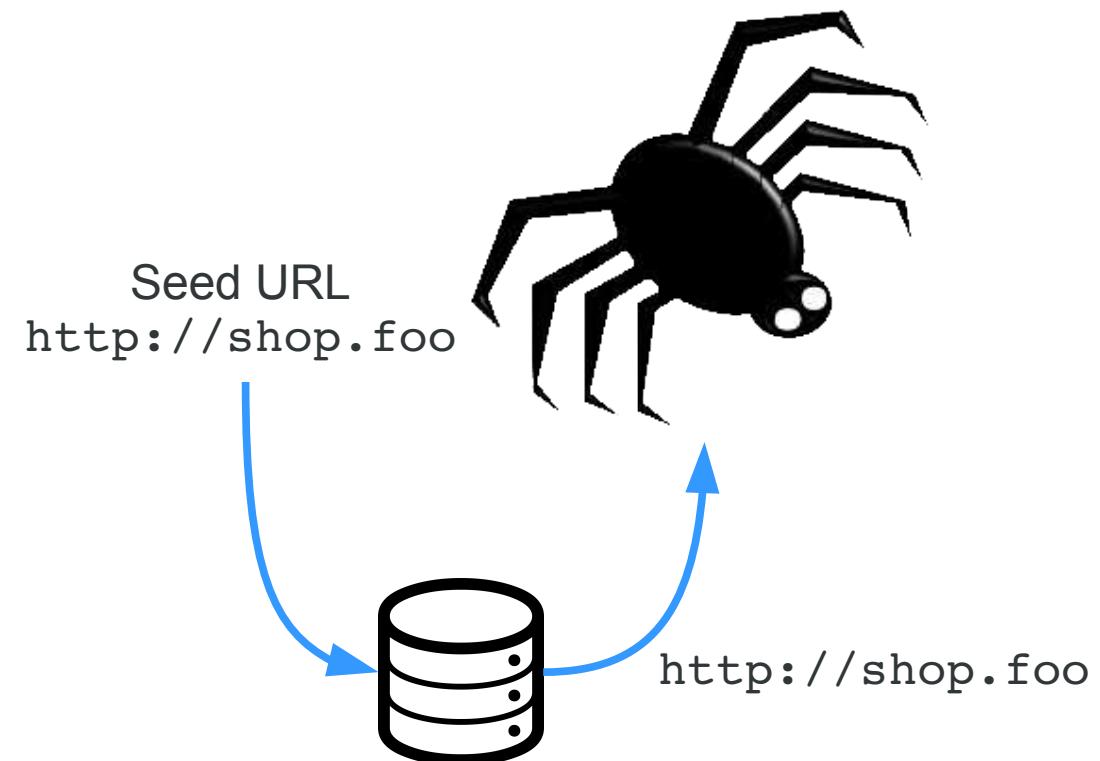


- (Semi-)automated security testing tools
- Follow a dynamic and black-box testing approach

Architecture

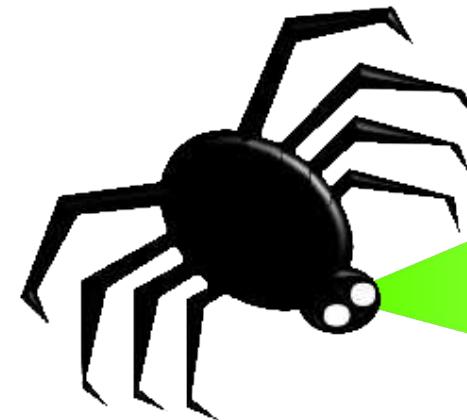


Crawler



Crawler

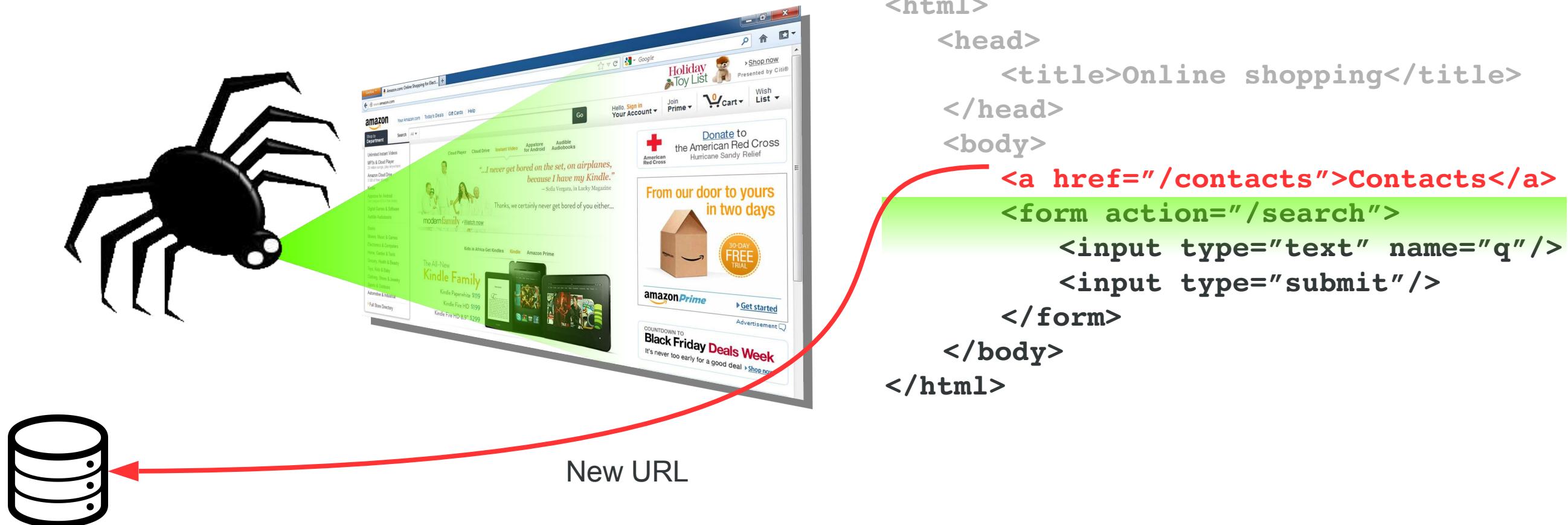
<http://shop.foo>



```
<html>
  <head>
    <title>Online shopping</title>
  </head>
  <body>
    <a href="/contacts">Contacts</a>
    <form action="/search">
      <input type="text" name="q"/>
      <input type="submit"/>
    </form>
  </body>
</html>
```

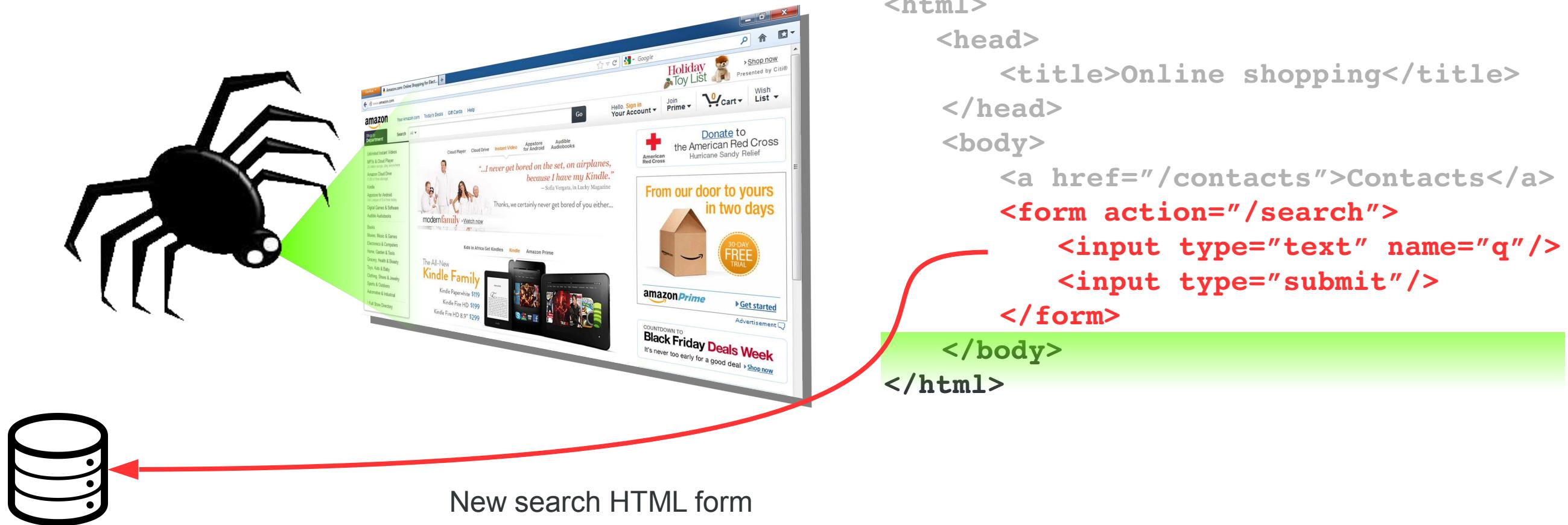
Crawler

http://shop.foo

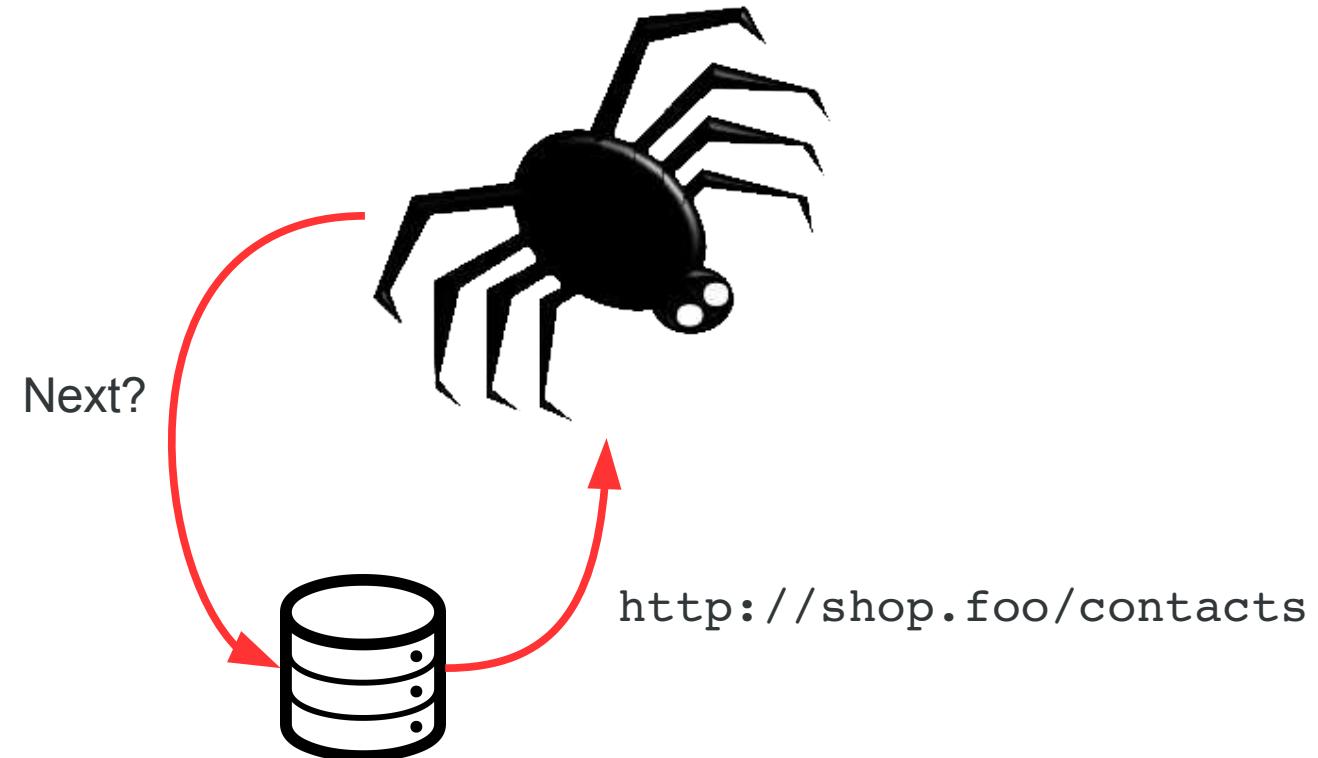


Crawler

http://shop.foo

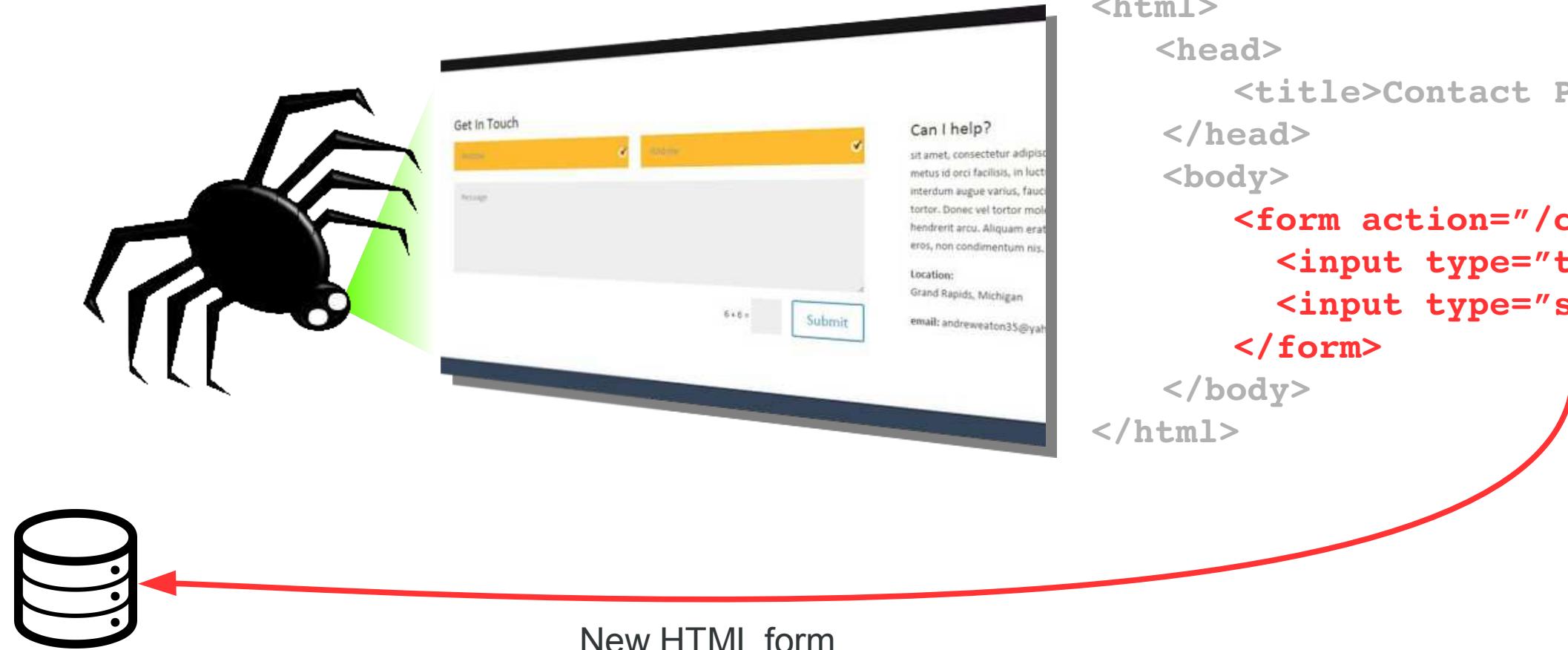


Crawler

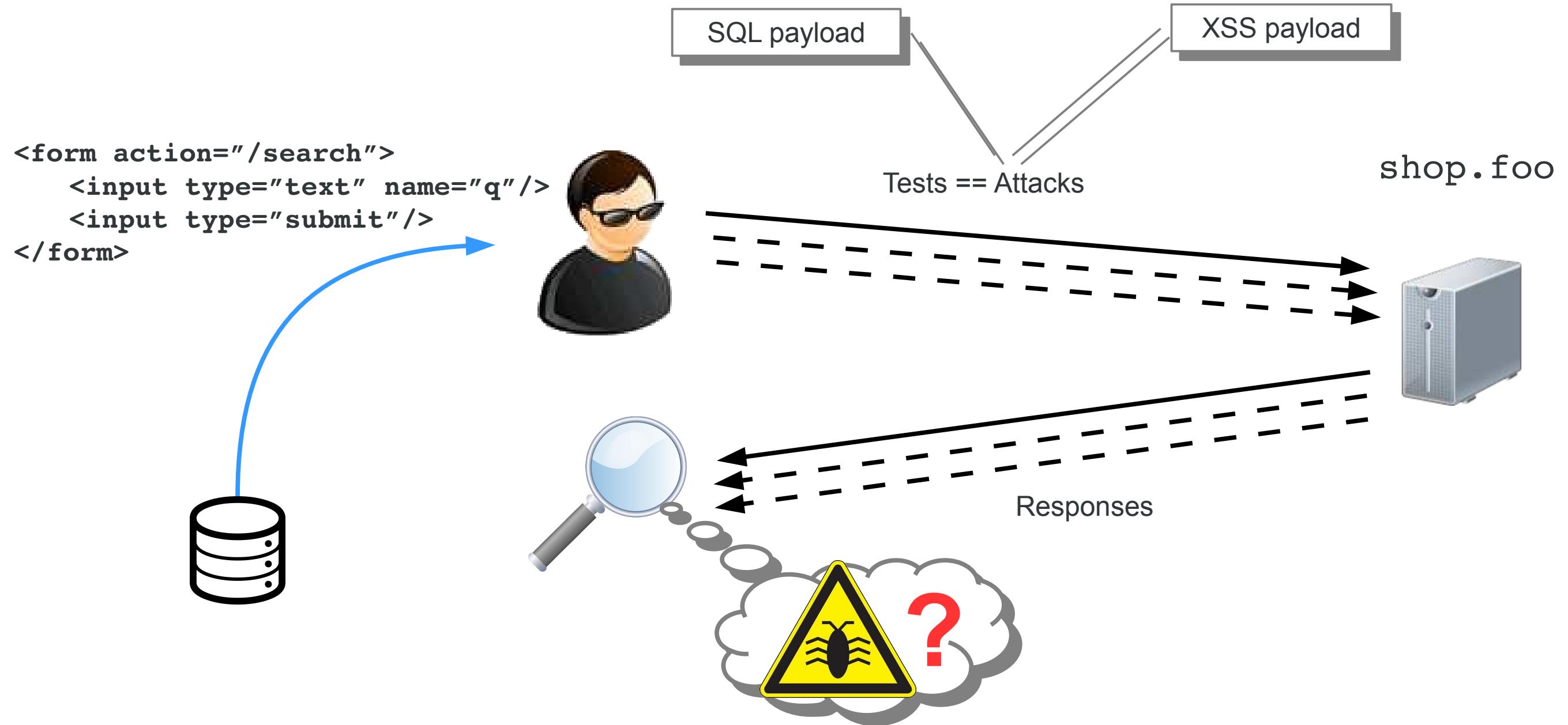


Crawler

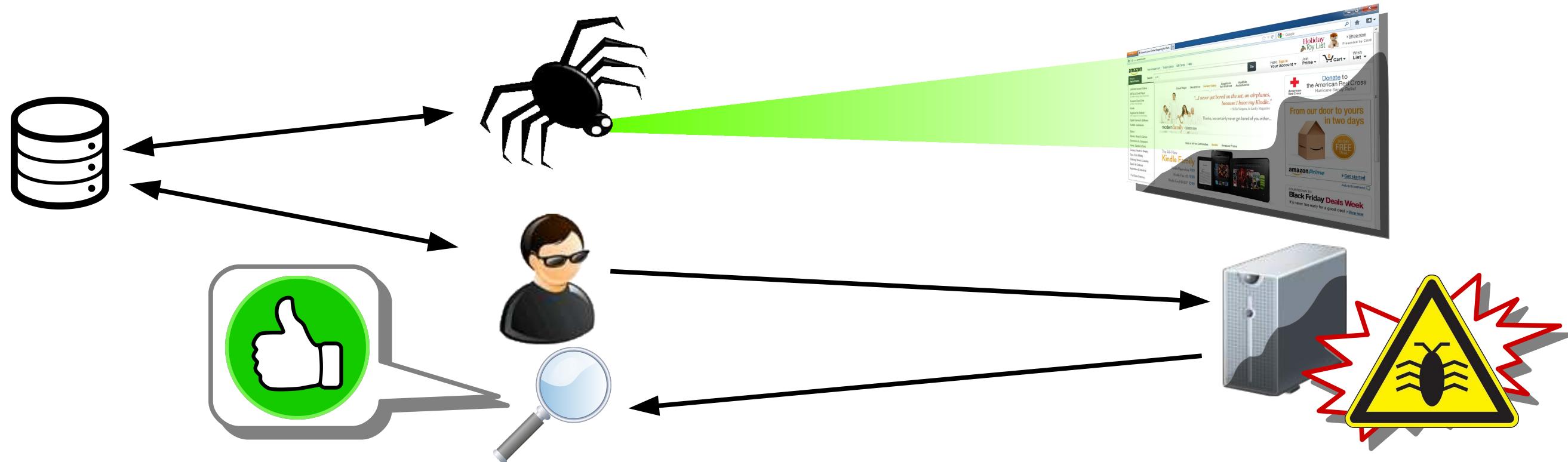
http://shop.foo/contacts



Security Testing



Crawler Critical for Coverage



- Crawler explores the Web application attack surface
 - Missing parts → missing possible vulnerabilities
- Existing crawlers based on:
 - HTML parsing and pattern matching to extract URLs
 - “clickable” areas to further explore the surface

Crawler and Modern Web Applications

- Complexity of client side has dramatically increased (i.e., stateful JS programs)

Crawler and Modern Web Applications

- Complexity of client side has dramatically increased (i.e., stateful JS programs)
- Links and forms can be built and inserted in the webpage at run-time

```
var url = scheme() + '://' + domain() + '/' + endpoint();
document.getElementById('myLink').href = url;
```

→ HTML parsing and pattern matching no longer sufficient

Crawler and Modern Web Applications

- Complexity of client side has dramatically increased (i.e., stateful JS programs)
- Links and forms can be built and inserted in the webpage at run-time

```
var url = scheme() + '://' + domain() + '/' + endpoint();
document.getElementById('myLink').href = url;
```

→ HTML parsing and pattern matching no longer sufficient

- JS is an event-driven language



- Functions executed upon events

→ Lack of support of event-based execution model

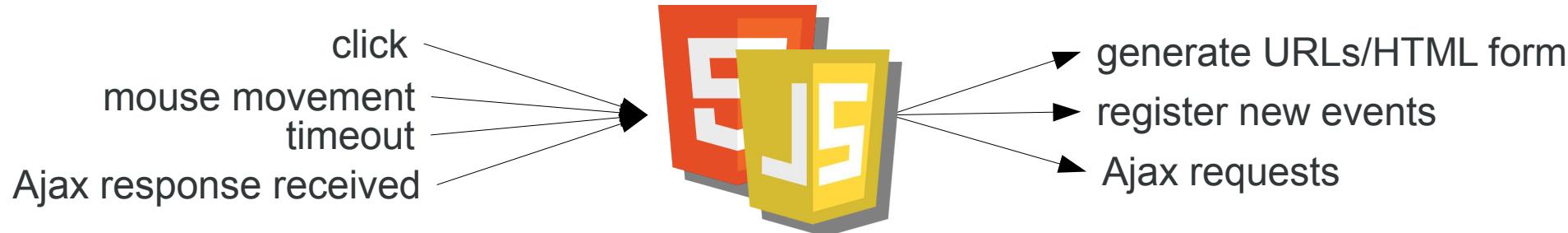
Crawler and Modern Web Applications

- Complexity of client side has dramatically increased (i.e., stateful JS programs)
- Links and forms can be built and inserted in the webpage at run-time

```
var url = scheme() + '://' + domain() + '/' + endpoint();
document.getElementById('myLink').href = url;
```

→ HTML parsing and pattern matching no longer sufficient

- JS is an event-driven language



- Functions executed upon events

→ Lack of support of event-based execution model

Large part of web applications remain unexplored!

Crawler and Modern Web Applications

- Complexity of client side has dramatically increased (i.e., stateful JS programs)
- Links and forms can be built and inserted in the webpage at run-time

```
var url = scheme() + '://' + domain() + '/' + endpoint();  
document.getElementById('myLink').href = url;
```

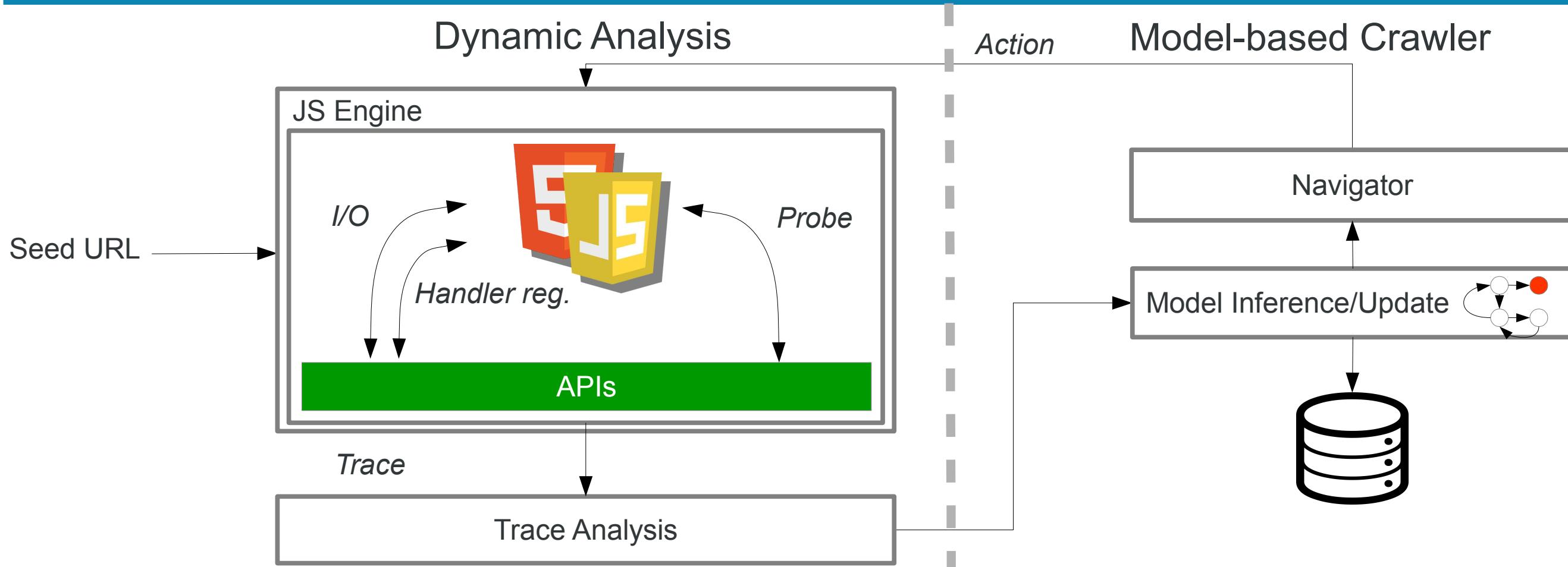
- We addressed the coverage problem with
 - JavaScript client side dynamic analysis
 - Model-based Crawler
- Build a tool: jÄk

• Functions executed upon events

→ Lack of support of event-based execution model

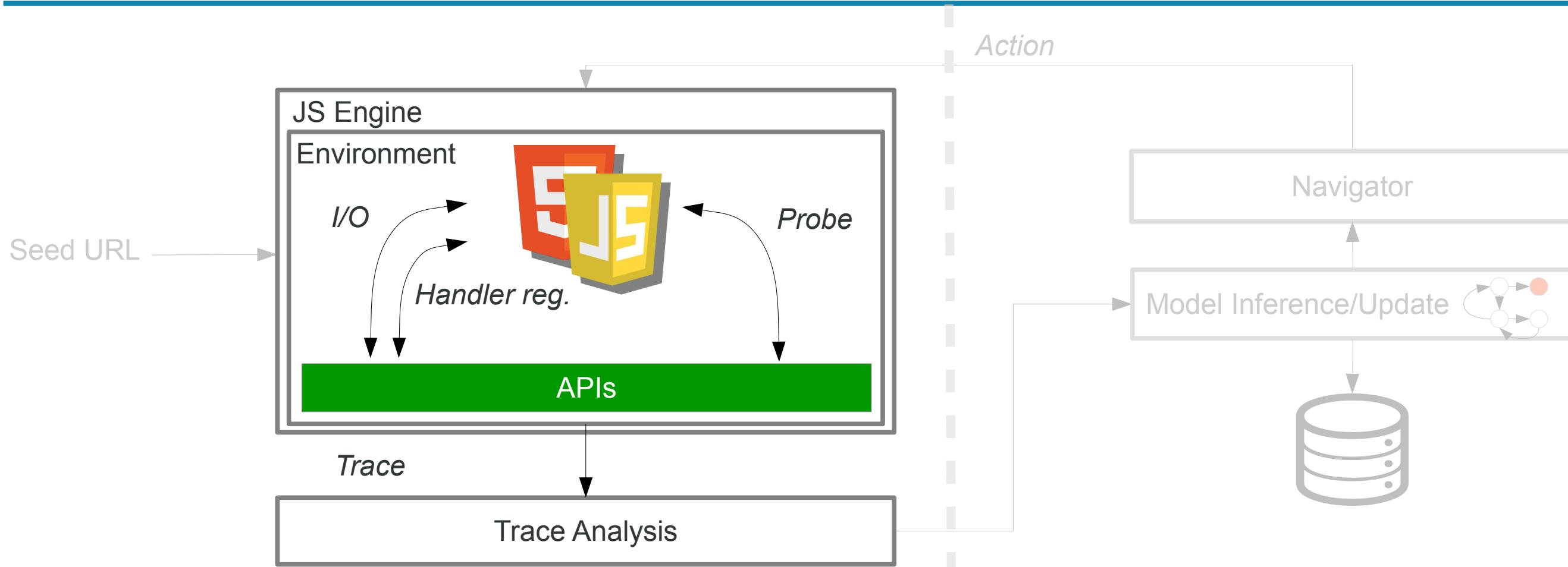
Large part of web applications remain unexplored!

Our Approach



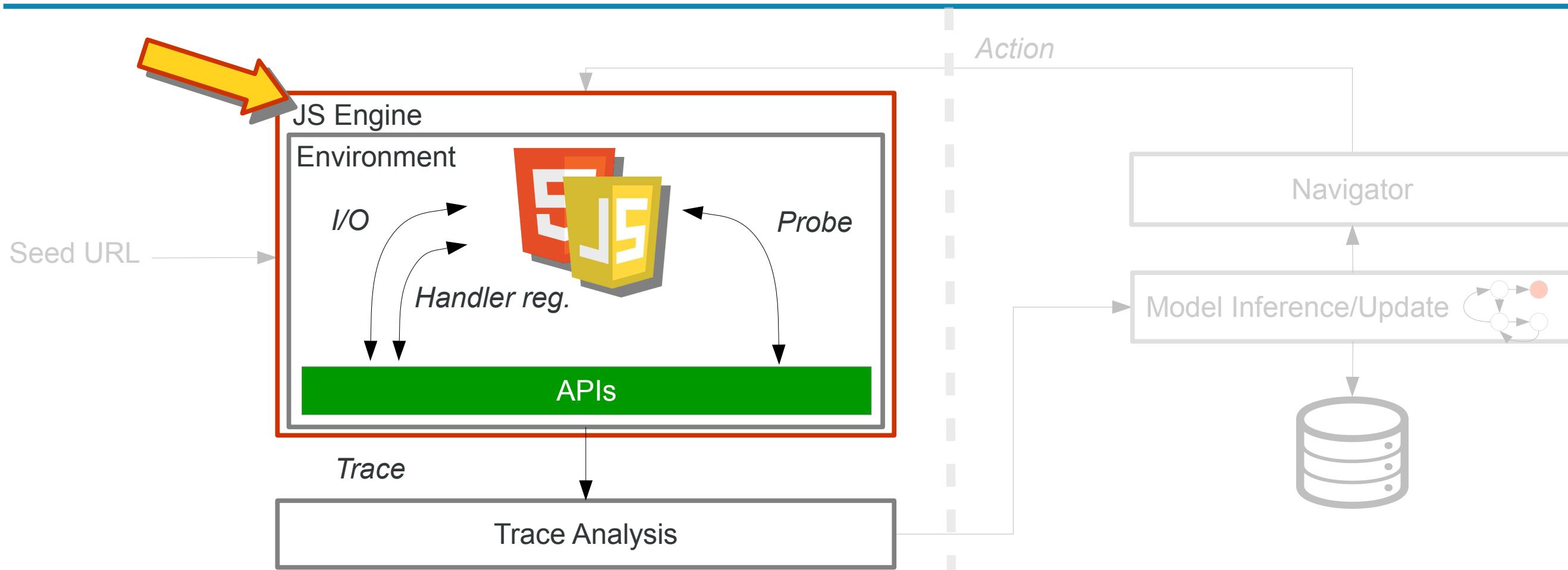
- Combine dynamic analysis with model-based crawler
 - Dynamic analysis monitors client side program execution
 - Crawler builds, maintains, uses a model of the visited attack surface

Dynamic Analysis



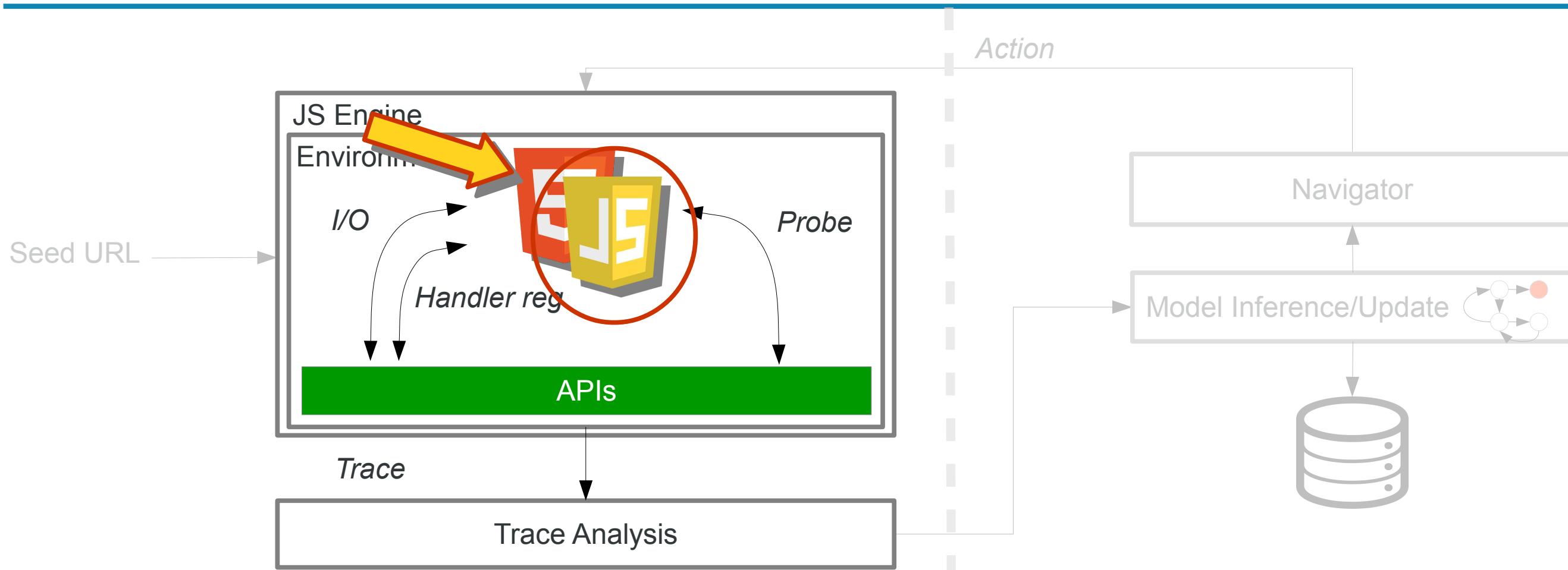
- Different approaches:

Dynamic Analysis



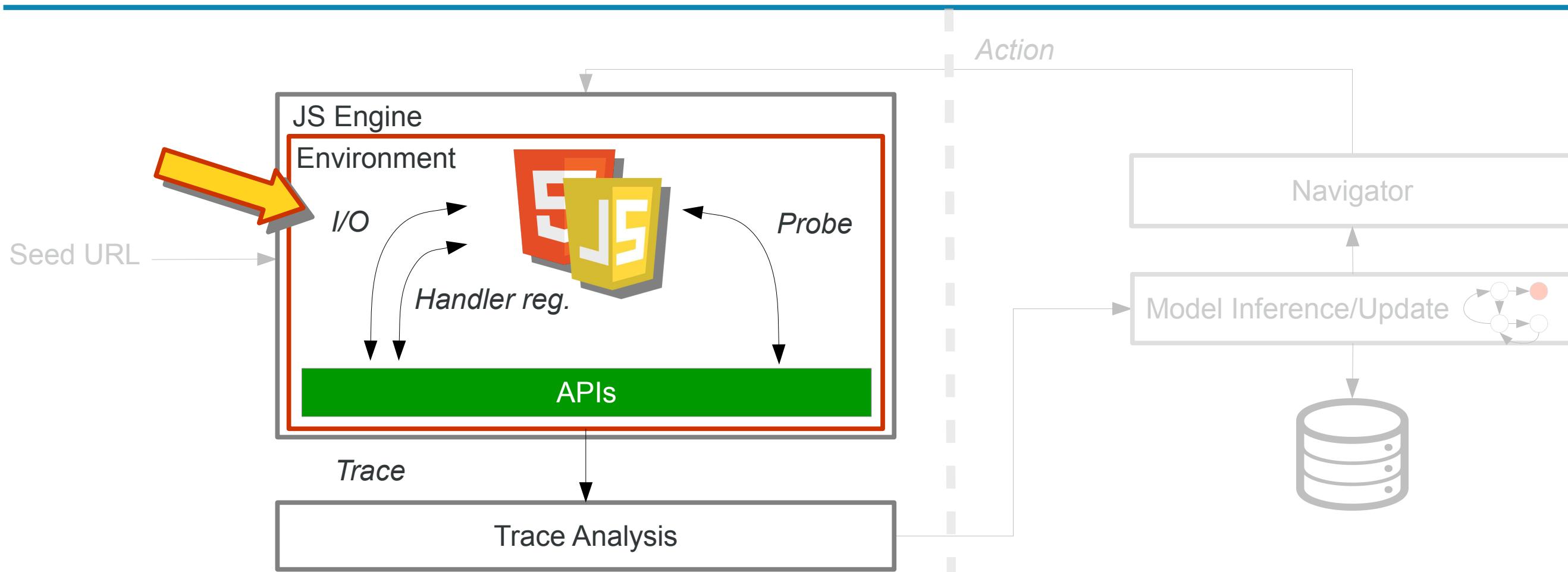
- Different approaches:
 - 1) JS engine instrumentation → laborious task, engine-dependent

Dynamic Analysis



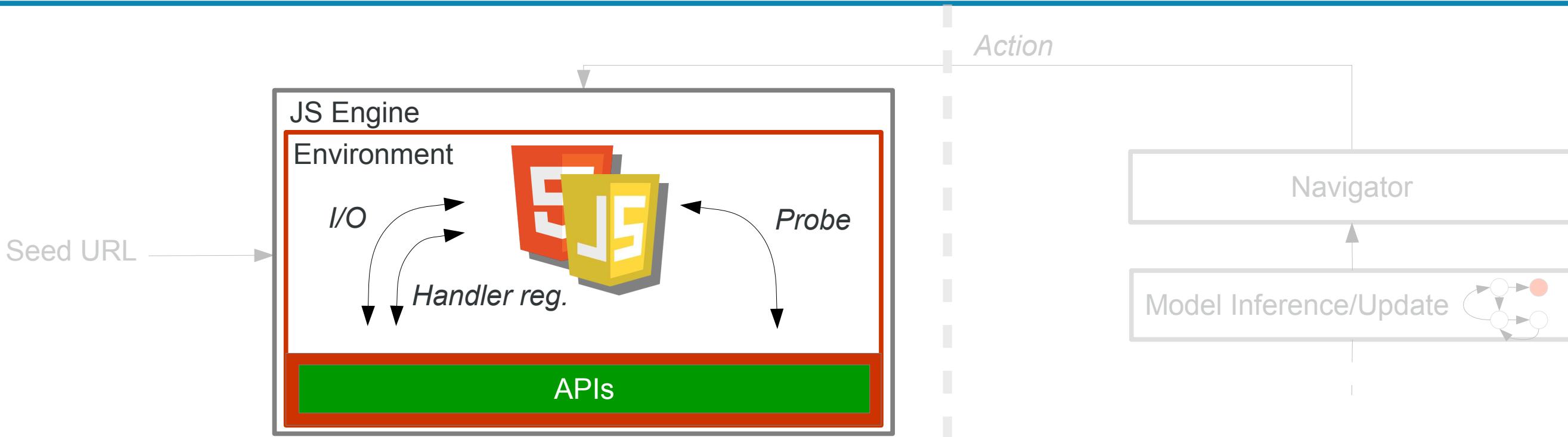
- Different approaches:
 - 1) JS engine instrumentation → laborious task, engine-dependent
 - 2) JS program instrumentation → JS code is not entirely available

Dynamic Analysis



- Different approaches:
 - 1) JS engine instrumentation → laborious task, engine-dependent
 - 2) JS program instrumentation → JS code is not entirely available
 - 3) Modification of execution environment

Dynamic Analysis



- Modify execution environment via **function hooking**:
 - Intercept API calls (e.g., network I/O and event handler registration)
 - Object manipulations (i.e., object properties)
 - Schedule DOM inspections
- Hooks installed by injecting own JS code:
 - Function redefinition
 - Set functions

Function Redefinition

Application
JS code

```
function handler() {  
    alert("hello world");  
}  
  
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



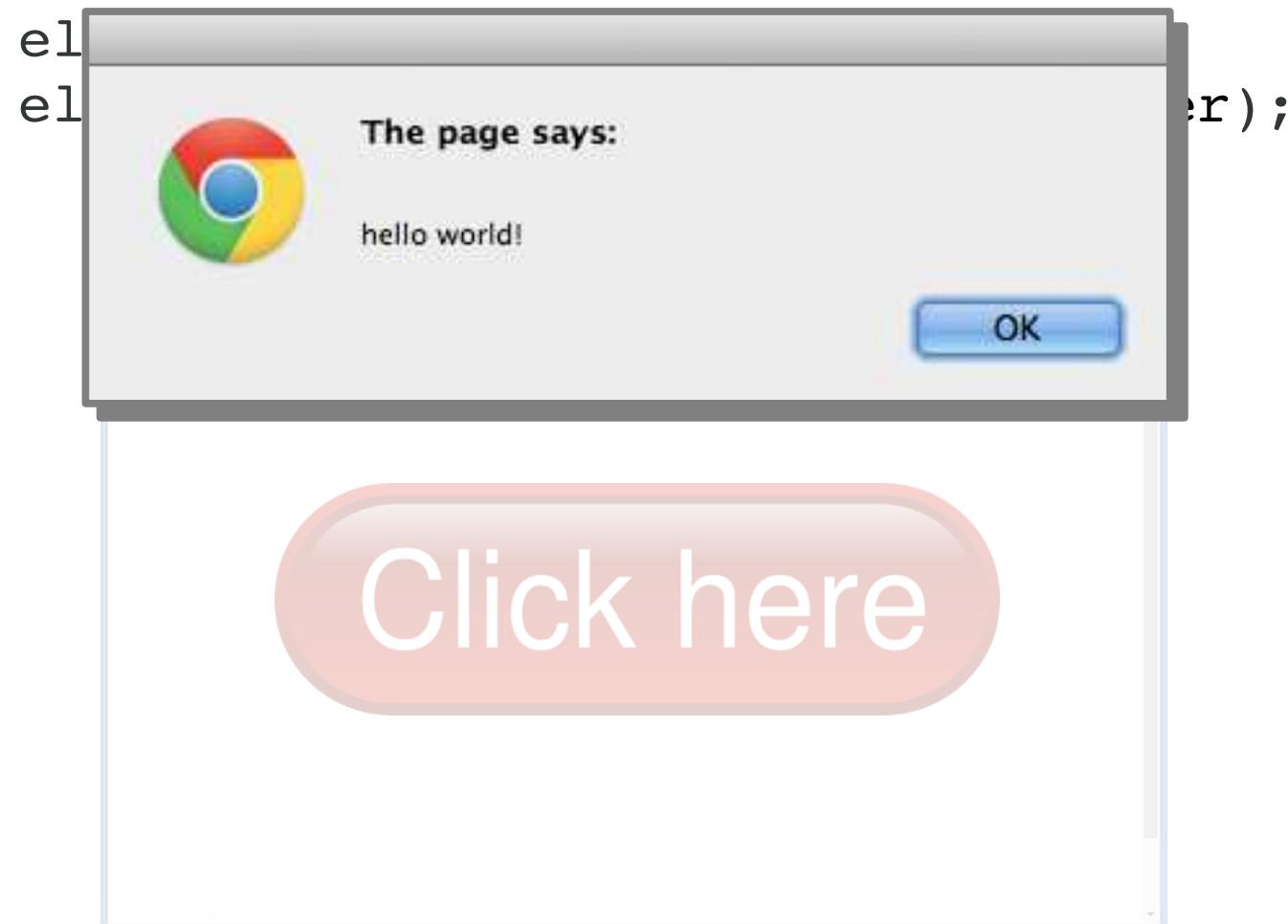
Function Redefinition

```
function handler() {  
    alert("hello world");  
}  
  
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



Function Redefinition

```
function handler() {  
    alert("hello world");  
}
```



Function Redefinition

```
function handler() {  
    alert("hello world");  
}
```

```
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



```
Element.prototype.addEventListener = function(e, h) {  
    [...]  
    listeners[e].append(h);  
}
```



Function Redefinition

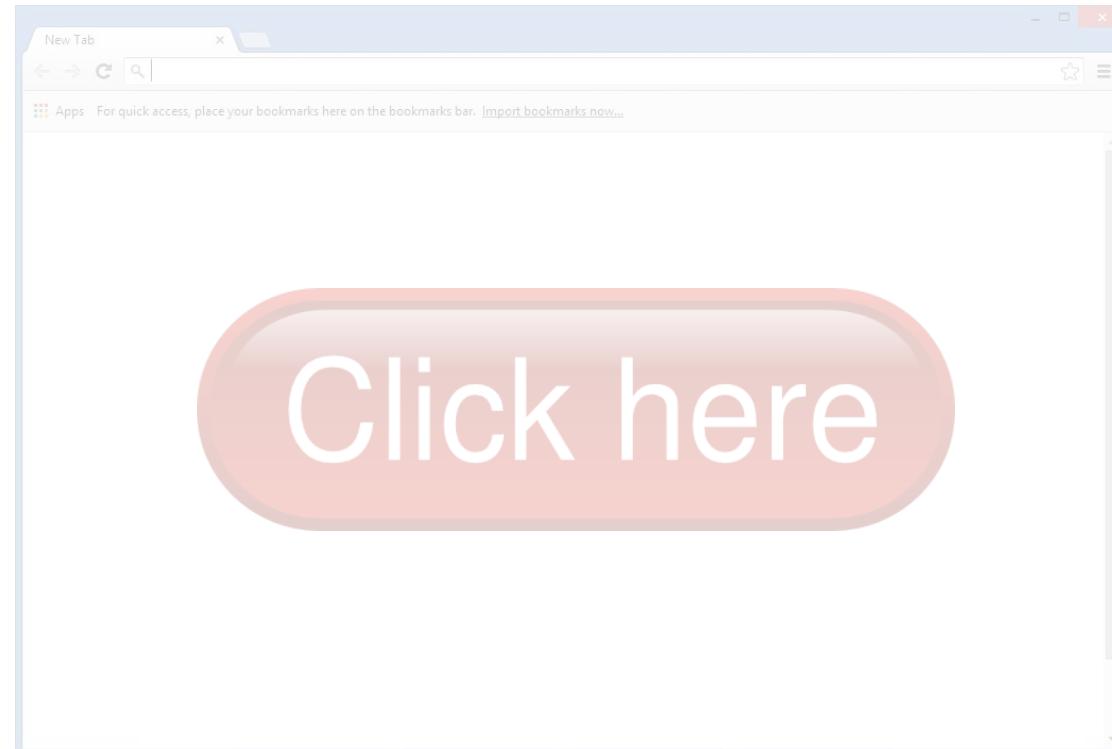
```
function handler() {  
    alert("hello world");  
}
```

```
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



```
Element.prototype.addEventListener = function(e, h) {  
    [...]  
    listeners[e].append(h);  
}
```

Intercept!

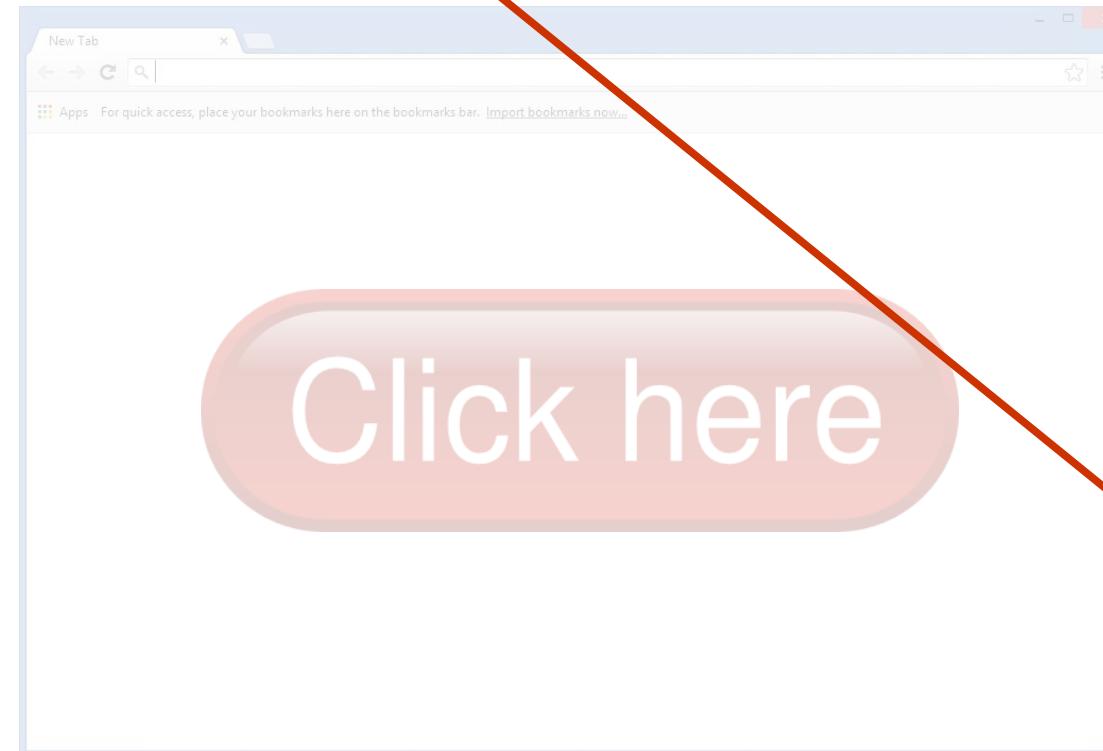


Function Redefinition

Application
JS code

```
preamble
function handler() {
    alert("hello world");
}

el = document.getElementById('img')
el.addEventListener("click", handler);
```



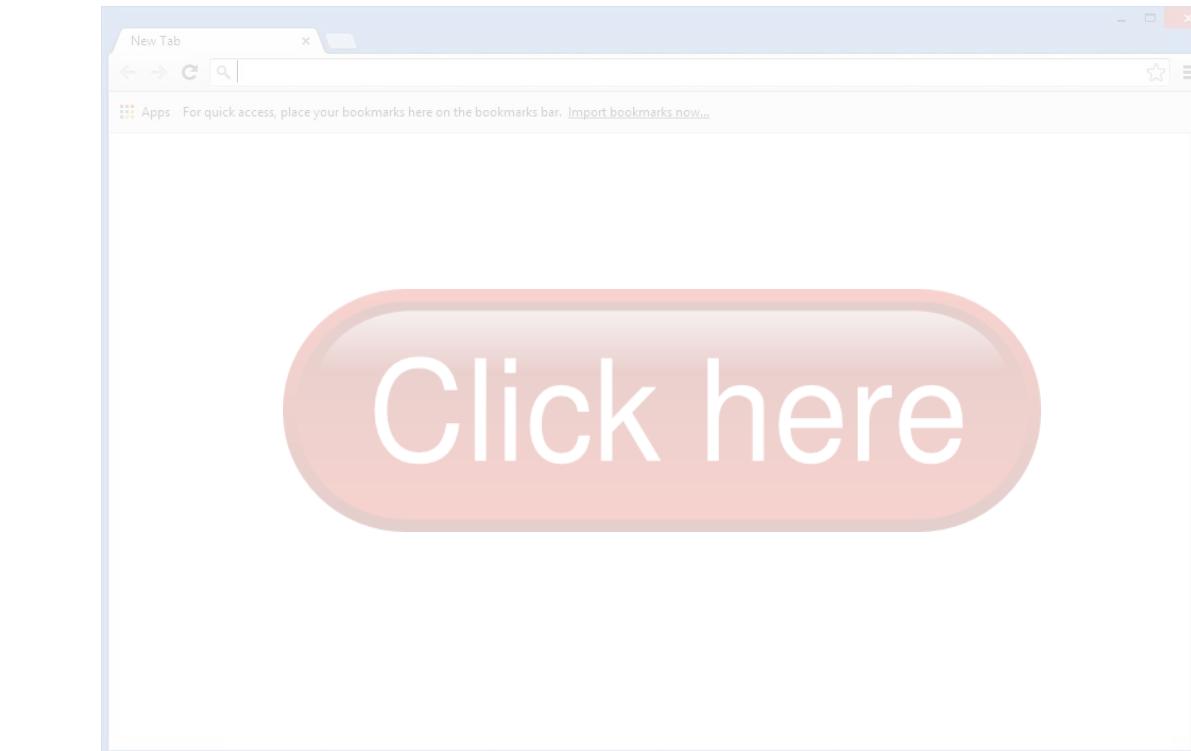
PREAMBLE

```
var orig_f = Element.prototype.addEventListener;
Element.prototype.addEventListener = function(){
    console.log("new handler registration");
    return orig_f.apply(this, argument);
};
```

Function Redefinition

preamble

```
function handler() {  
    alert("hello world");  
}  
  
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



API

```
Element.prototype.addEventListener = function(e, h) {  
    [...]  
    listeners[e].append(h);  
}
```

PREAMBLE

```
var orig_f = Element.prototype.addEventListener;  
  
Element.prototype.addEventListener = function(){  
    console.log("new handler registration");  
    return orig_f.apply(this, argument);  
};
```

Function Redefinition

preamble

```
function handler() {  
    alert("hello world");  
}  
  
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



API

```
Element.prototype.addEventListener = function(e, h) {  
    [...]  
    listeners[e].append(h);  
}
```

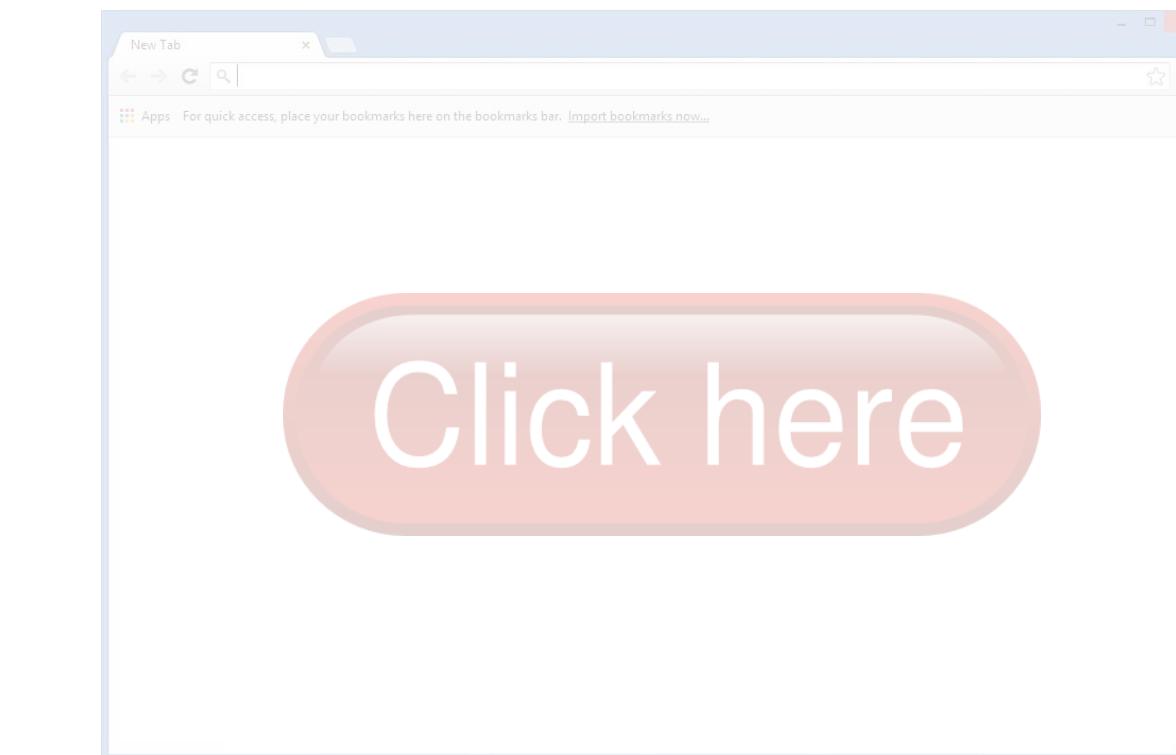
PREAMBLE

```
var orig_f = Element.prototype.addEventListener;  
  
Element.prototype.addEventListener = function(){  
    console.log("new handler registration");  
    return orig_f.apply(this, argument);  
};
```

Function Redefinition

preamble

```
function handler() {  
    alert("hello world");  
}  
  
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



API

```
Element.prototype.addEventListener = function(e, h) {  
    [...]  
    listeners[e].append(h);  
}
```

PREAMBLE

```
var orig_f = Element.prototype.addEventListener;  
  
Element.prototype.addEventListener = function(){  
    console.log("new handler registration");  
    return orig_f.apply(this, argument);  
};
```

Function Redefinition

preamble

```
function handler() {  
    alert("hello world");  
}  
  
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



API

```
Element.prototype.addEventListener = function(e, h) {  
    [...]  
    listeners[e].append(h);  
}
```

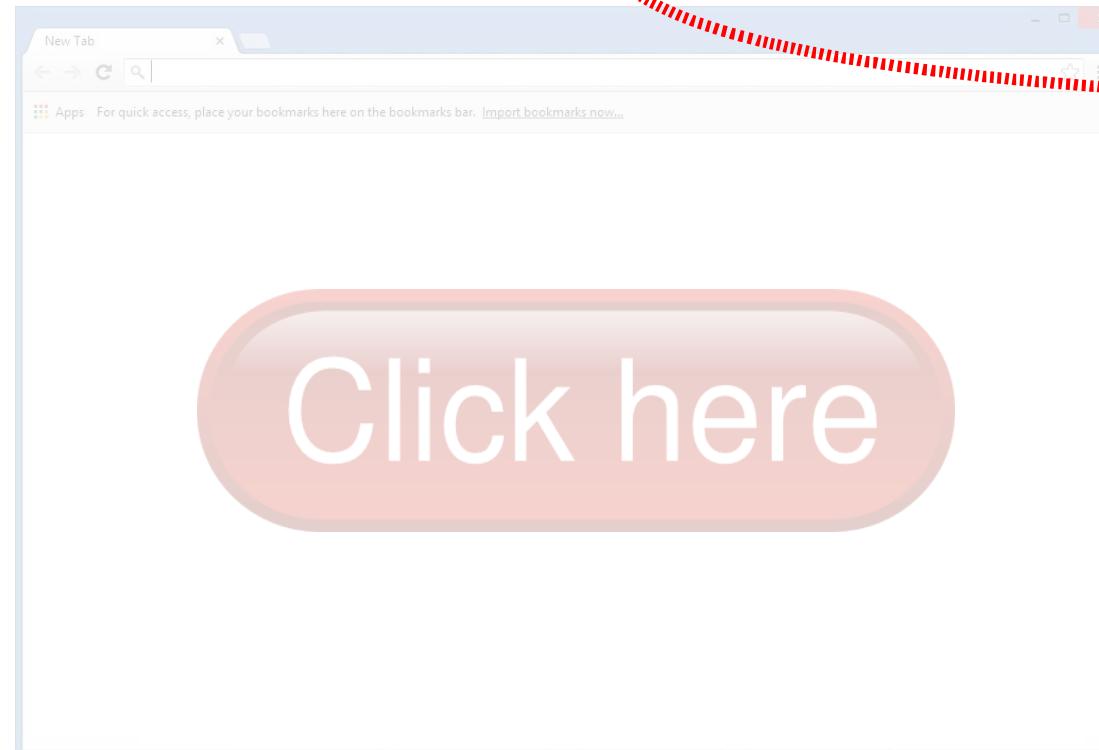
PREAMBLE

```
var orig_f = Element.prototype.addEventListener;  
  
Element.prototype.addEventListener = function(){  
    console.log("new handler registration");  
    return orig_f.apply(this, argument);  
};
```

Function Redefinition

preamble

```
function handler() {  
    alert("hello world");  
}  
  
el = document.getElementById('img')  
el.addEventListener("click", handler);
```



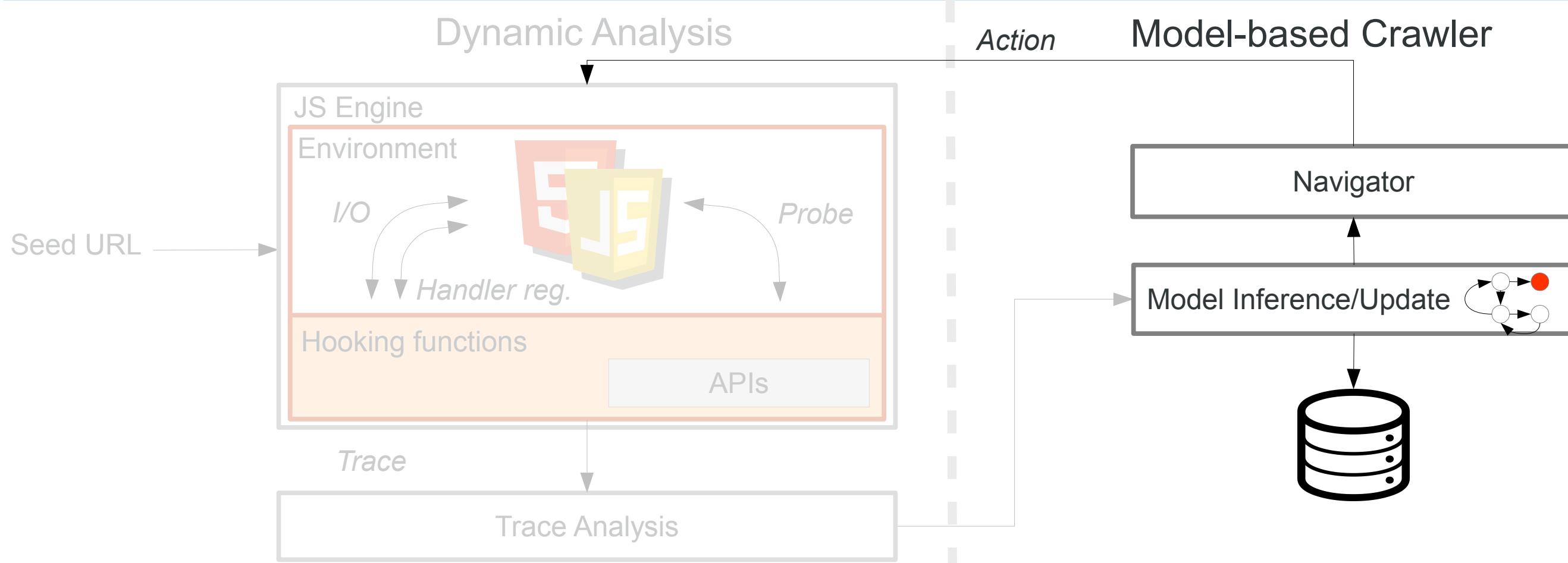
API

```
Element.prototype.addEventListener = function(e, h) {  
    [...]  
    listeners[e].append(h);  
}
```

PREAMBLE

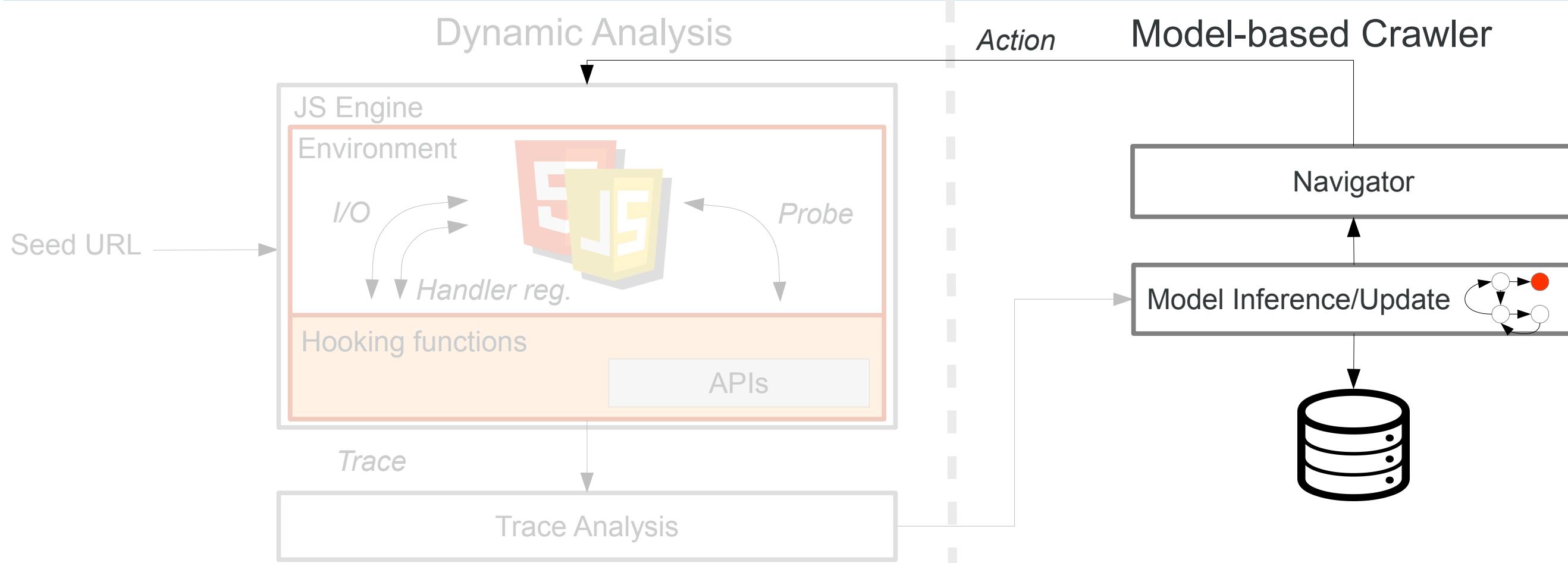
```
var orig_f = Element.prototype.addEventListener;  
  
Element.prototype.addEventListener = function(){  
    console.log("new handler registration");  
    return orig_f.apply(this, argument);  
};
```

Model-based Crawler



- Creates and maintain a web application model
 - Oriented graph: nodes are page clusters and edges are URLs, HTML forms, or events
- Model used to decide on the next action
 - Priority: Events → high, URLs/forms → low

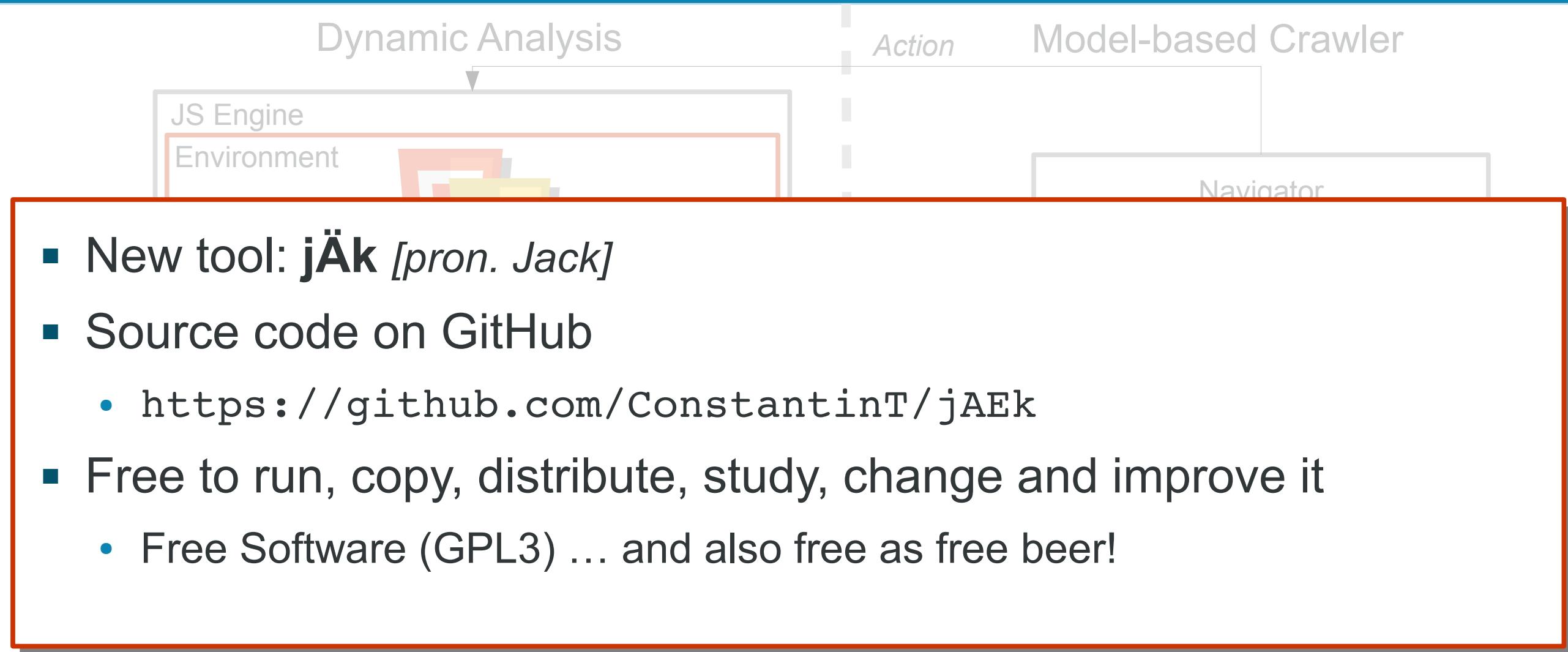
Model-based Crawler



- Creates and maintain a web application model
 - Oriented graph: nodes are page clusters and edges are URLs, HTML forms, or events
- Model used to decide on the next action
 - Priority: Events → high, URLs/forms → low

Assessment

Our Tool

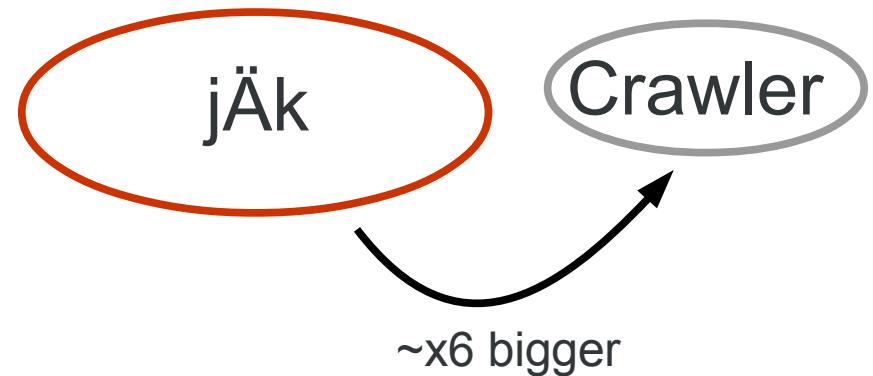


Experiments

- Comparative analysis
 - Skipfish, W3af, Wget, and Crawljax
- Case studies:
 - WIVET (Web Input Vector Extractor Teaser)
 - assess strength and limitations of existing crawlers
 - 13 web applications
 - studied coverage and vulnerability detection power

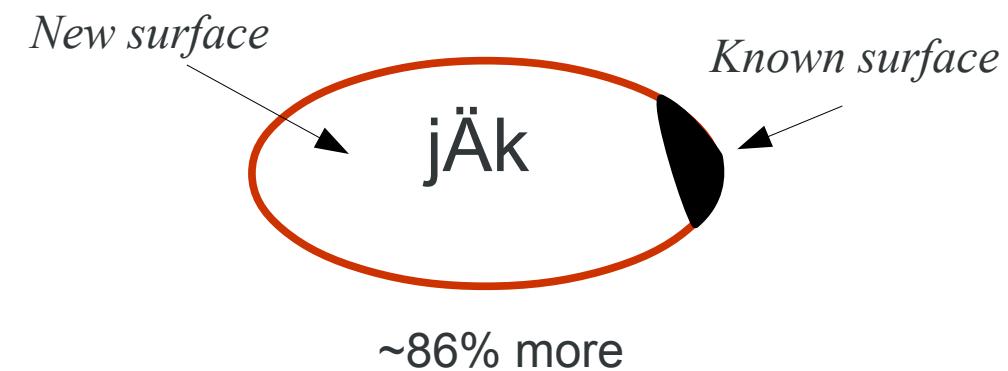
Coverage

- Explored surface by jÄk (# of unique URL structs.)
 - x16 (Crawljax) to x2 (Skipfish) bigger



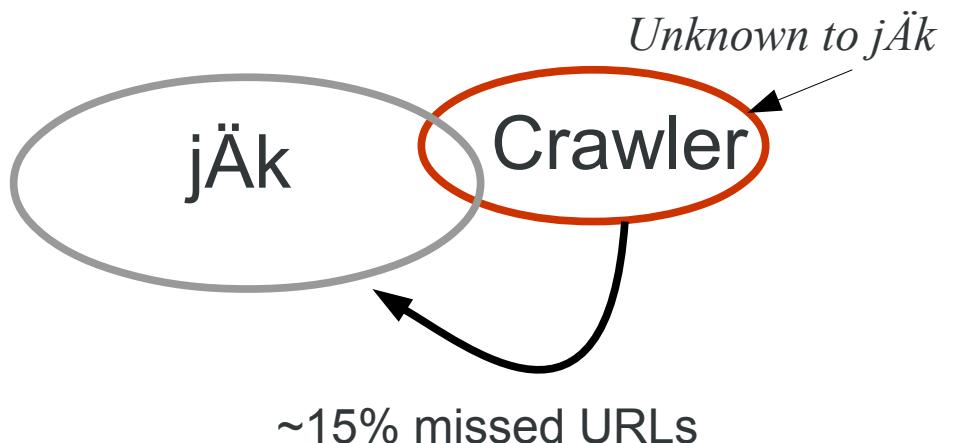
Coverage

- Explored surface by jÄk (# of unique URL structs.)
 - x16 (Crawljax) to x2 (Skipfish) bigger
- Relative size of new surface:
 - From +70% (Wget) to +98% (Crawljax) of URLs are new



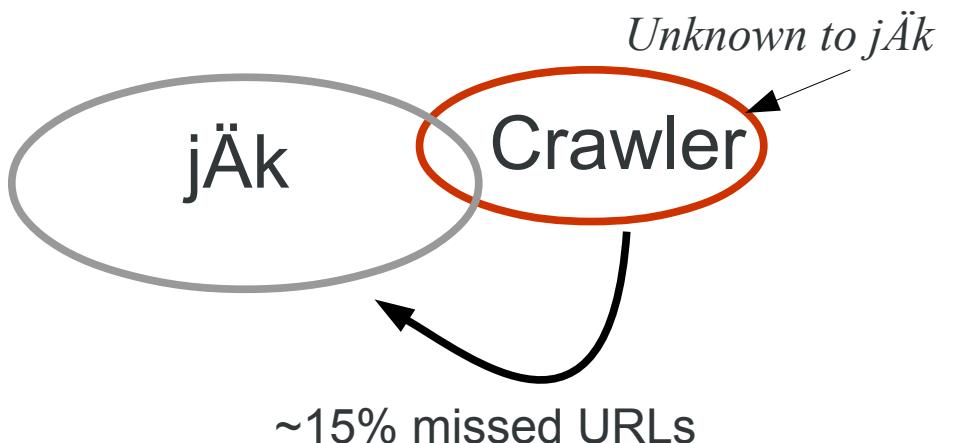
Coverage

- Explored surface by jÄk (# of unique URL structs.)
 - x16 (Crawljax) to x2 (Skipfish) bigger
- Relative size of new surface:
 - From +70% (Wget) to +98% (Crawljax) of URLs are new
- Global surface missed by jÄk:
 - From 22% (Skipfish) to 0.5% (Crawljax) are missed



Coverage

- Explored surface by jÄk (# of unique URL structs.)
 - x16 (Crawljax) to x2 (Skipfish) bigger
- Relative size of new surface:
 - From +70% (Wget) to +98% (Crawljax) of URLs are new
- Global surface missed by jÄk:
 - From 22% (Skipfish) to 0.5% (Crawljax) are missed
 - Further analysis:
 - 75% of missed are due to URL forgery
 - 25% to static resources, unsupported action, and others



Conclusion

Conclusion/Takeaway

- Novel technique based on
 - **dynamic analysis** of JS program + **model-based crawling**
- Built **jÄk**, a tool implementing our approach
- Assessed against 13 web applications
- Our result show that jÄk explores a surface ~6x larger with +86% new URLs