



# Wikispam, Wikispam, Wikispam PmWiki

---

Patrick R. Michaud, Ph.D.

March 4, 2005

[Google Search](#)

[I'm Feeling Oppressed](#)



# Talk Outline

---

- What is wikispam?
- Know thy opponent(s)
- Countermeasures
- PmWiki security options

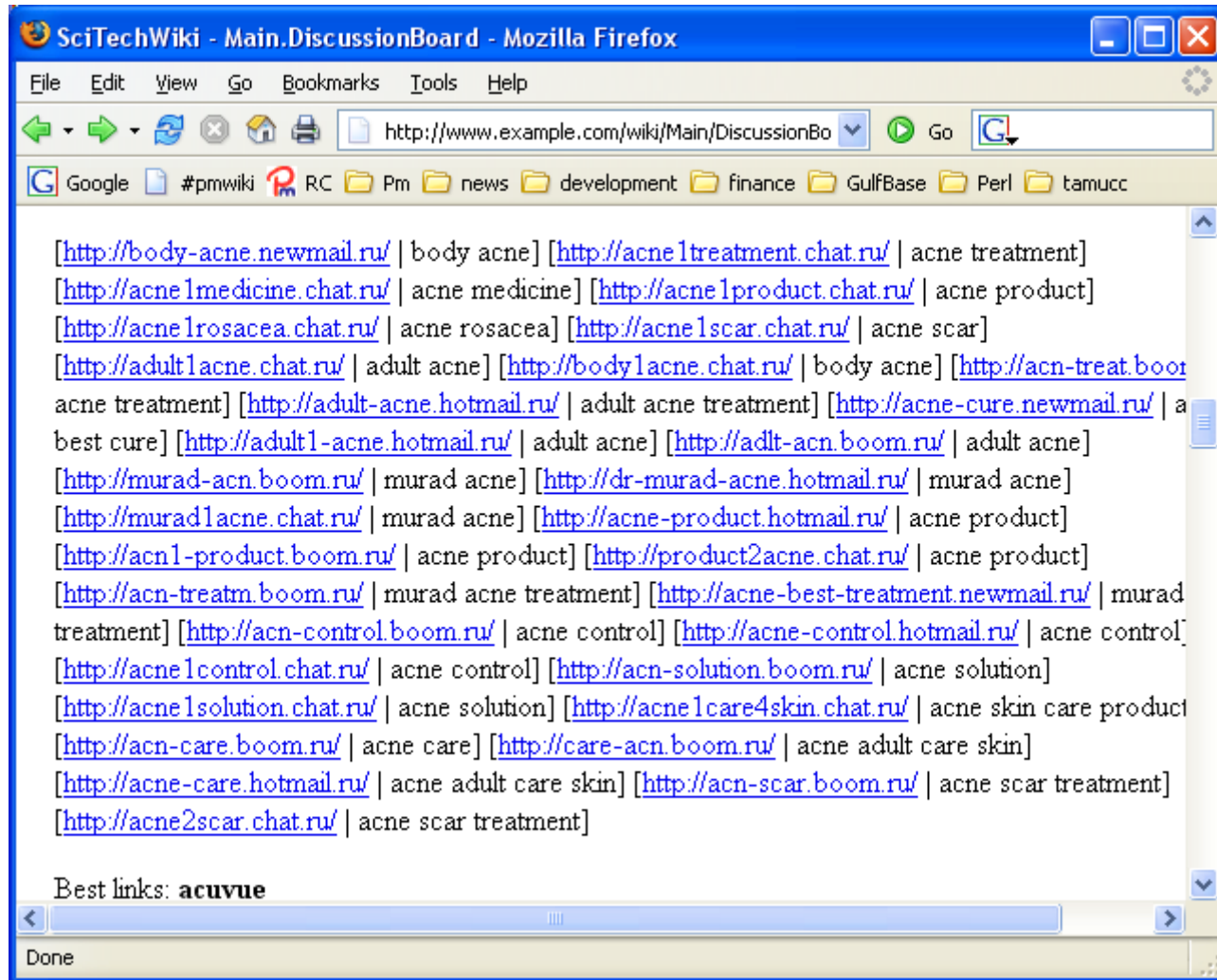


# What is Wikispam?

---

- Wikispam, or "comment spam", refers to unwanted external links in publicly writable web pages
  - Wikis
  - Blogs
  - Other

# What is Wikispam?



# Know thy opponent(s) - 1

---

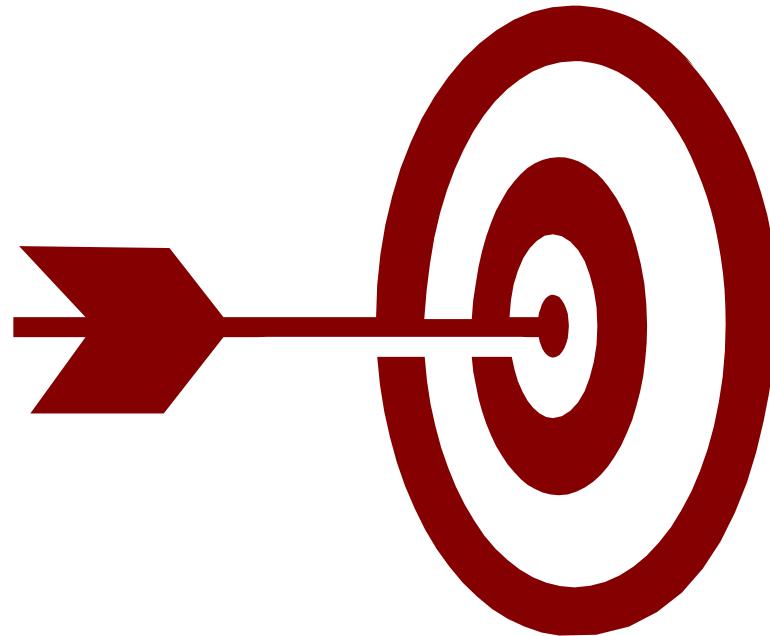
- Why do people spam?
  - Vandalism/malice
  - Money

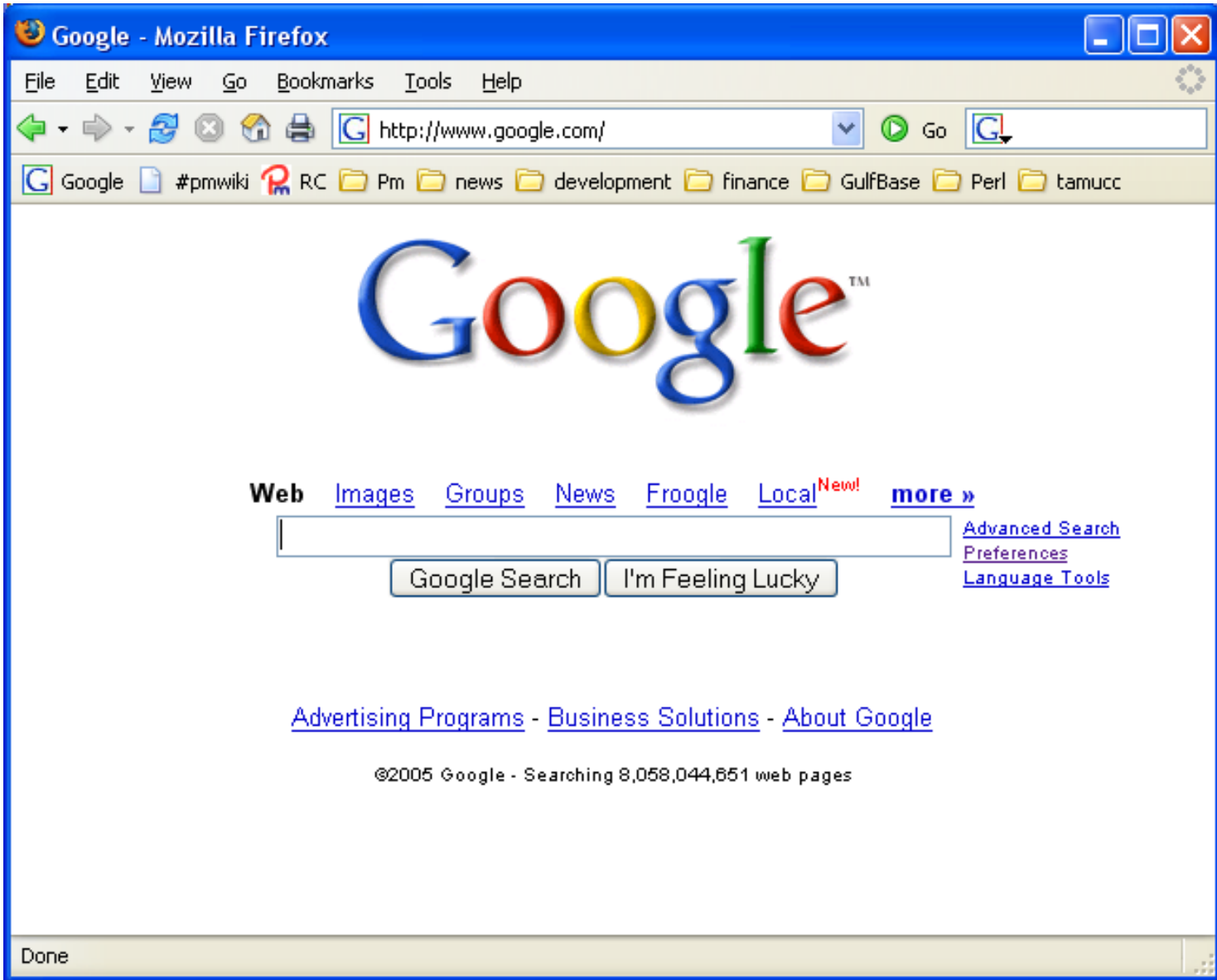




# Wikispam's primary target

---





Google - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://www.google.com/ Go

Google #pmwiki RC Pm news development finance GulfBase Perl tamucc



Web Images Groups News Froogle Local<sup>New!</sup> more »

Google Search I'm Feeling Lucky

Advanced Search Preferences Language Tools

Advertising Programs - Business Solutions - About Google

©2005 Google - Searching 8,058,044,851 web pages

Done



# The Googlephant in the room

---

- Companies and organizations want high Google rankings
- Google listings are ordered according to "PageRank"
  - Each link from page A to page B is counted as a "vote for B"
  - In general, the more pages that vote for B, the higher B's rank
- So, for a high page rank, get lots of pages to link to your site, then Google will list you first



There's names for this...

---

GoogleBombs

"Search Engine Optimization"



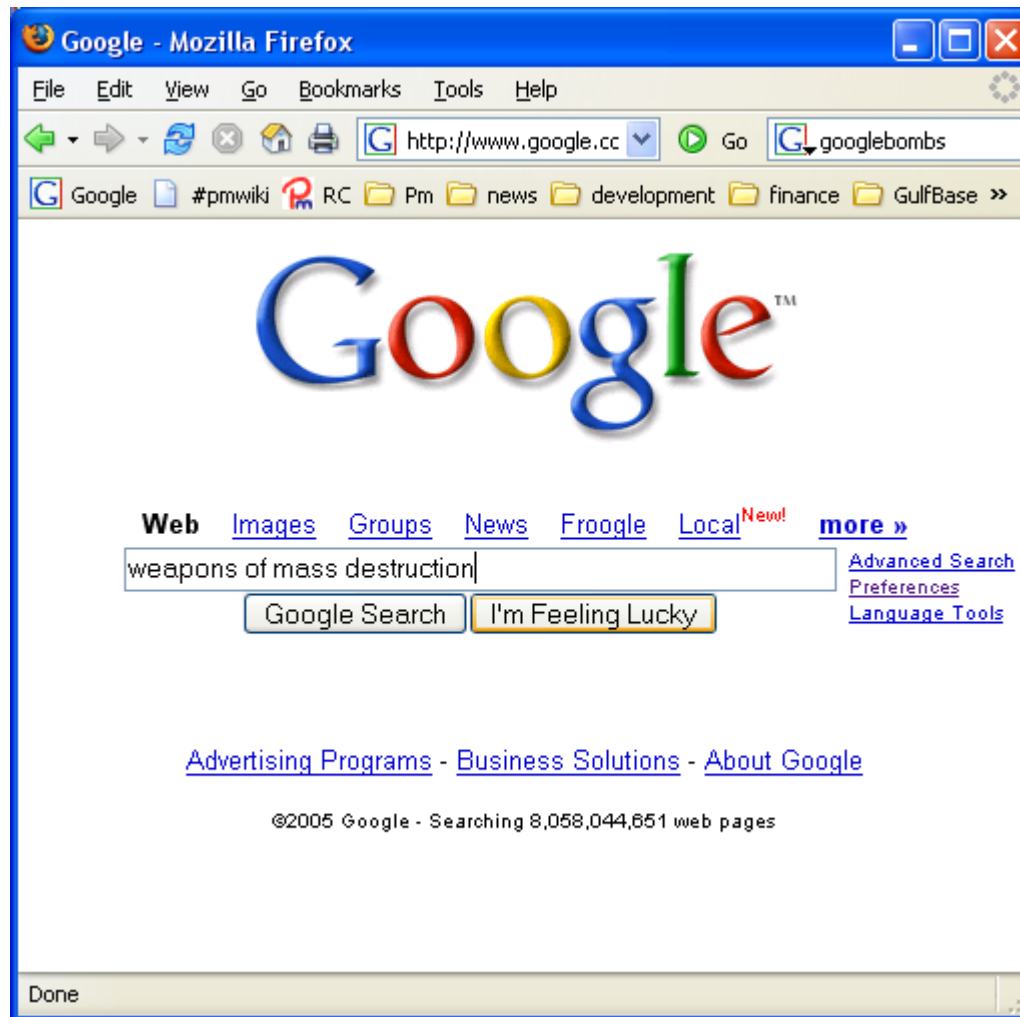
# Make\$Money\$Fast

---

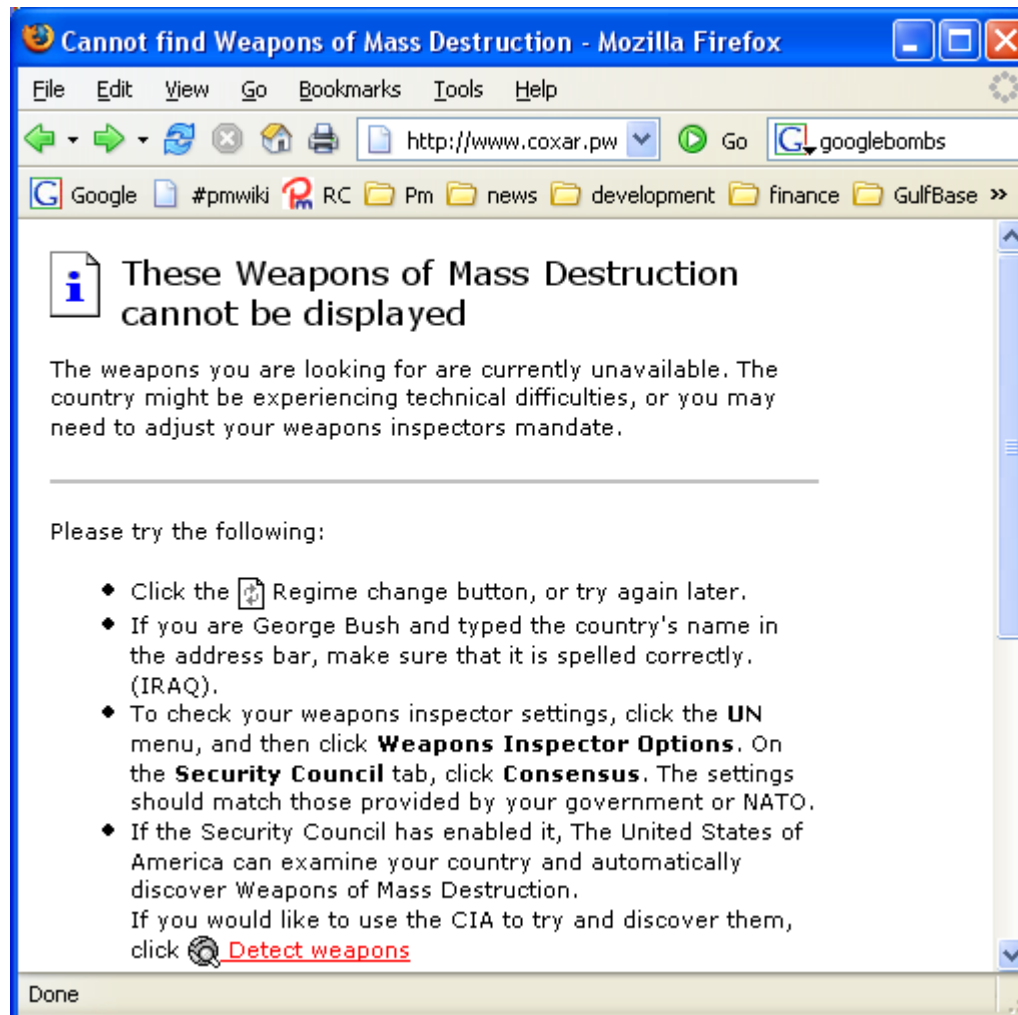
So, what's the easiest way to generate lots of links to a page...?

Add them to blogs, wikis, anywhere that has unsecured web content!

# Links of Mass Deception



# Links of Mass Deception





# Finding spam

---

- Wikispammers tend to flood pages with lots of links
- They may attempt to hide links in unexpected places (periods and punctuation)
- Writing a bot to generate such links is easy



# What's a wiki to do?

---

- Countermeasures
  - Get rid of that \$!#%\* wiki!
  - Open but vigilant (a.k.a. "live with it")
  - Blocklists
  - Spam filters
  - Captchas
  - Access controls
  - URL approvals
  - Learning to love Google all over again



# Countermeasure 1 - Elimination

---

aka "Whose stupid idea was this, anyway?"

- Pros: No more wikispam
- Cons: No more wiki



# Countermeasure 2 - Openness

---

aka "Maybe we can outrun the bots!"

- Rely on openness of wiki to enable lots of people to correct spam when its found
- Pro: wiki makes it simple to undo spam when it's found
- Con: people get tired of cleaning up other's messes
- Wikipedia takes this approach



## Countermeasure 3 - Blocklists

---

aka "They don't work for email,  
might as well try them here!"

- Keep track of spammers IP address, then block posts from that IP address
- Er, IP address *range*, that is
- Pro: simple to implement
- Cons:
  - Maintaining lists of IPs to block
  - Sometimes block legitimate IPs
  - Spammers know how to jump IPs



# Countermeasure 4 - Spamfilters

---

aka "What doesn't work for email  
is good enough for us!"

- Filter based on keywords and page analysis
- Bayesian networks
- Pros: could be effective
- Cons: never is effective
  - Too much effort to train filters
  - Hard to classify spam/non-spam
  - Spammers know how to adapt to filters

# Countermeasure 5 - Captchas

"The bot can't hurt what it can't see!"

- Captcha - telling computers and humans apart automatically



- Theory: image (and character) recognition is hard for computers, easy for humans
- Thus, a simple test can distinguish the two and filter bots
- Used by many commercial sites



# Captchas may not work

---

- Inconvenient for editing - new captcha required for each post
- Bots can solve captchas
  - Captcha is much simpler than OCR
  - Many captcha systems have been broken
- Spammers can enlist humans to solve captchas on their behalf
- Some humans cannot solve them
- Reduces web accessibility



GG8a



# Countermeasure 6 - access ctrls

---

When "open" is just too much of a good thing.

- Add passwords or other access controls to wiki pages
- Pros: Well known, works well
- Cons:
  - Increased administrative overhead
  - Loses the "open nature" of the wiki
  - Annoying



# Countermeasure 7 - url approvals

---

aka "Getting to the source of the problem"

- Block spammers ability to create links
- Known domains and urls are "approved"
- Any urls not in the approved list are not converted to links
- Provide an easy mechanism to approve urls



# Url approvals

---

- Pros:

- denies spammers reward for effort
- allows open editing

- Cons:

- doesn't prevent spam -- still has to be cleaned up
- maintenance of "approved urls" list



# Url approvals + blocking

---

"Prevention is the best defense."

- Wikispammers tend to post lots of unapproved links
- Authors generally have only a few
- Solution: Block posts that have more than a given number of unapproved external links
- Pros: Prevents spam from appearing
- Cons: May temporarily inconvenience authors wanting to post lots of links
- Since implementing on pmwiki.org in Dec 2004, spam hasn't been a problem



# Learning to love Google again

---

- Google and other search engine companies recognized that they were a big part of the problem
- Decided to become part of the solution
- Sites can designate links that should not count for search-engine ranking with  
`<a href='...' rel='nofollow'>`



`<a href='...' rel='nofollow'>`

---

- Despite the term 'nofollow', only means that search engines do not count the link as a vote (the links can still be followed)
- Used only on links to external sites where the validity of the link is unknown
- Pros: cuts benefit of spamming off at source
- Cons:
  - may take some time for spammers to notice (or care)
  - takes time for sites to upgrade



```
<a href='...' rel='nofollow'>
```

---

- PmWiki and other engines have quickly adopted this convention
- PmWiki allows rel='nofollow' only on unapproved links



# More about PmWiki security

---