

Effective Email Spam Control from Traffic Analysis

Richard Clayton

JANET NETWORKSHOP 37
Cambridge, 1st April 2009



A talk about ISP mail handling

BUT, this audience not all that different!

- Outgoing log processing
 - spot problems on your smarthost
- Incoming log processing
 - spot email being sent “direct”
- Aardvarks & Zebras
 - different people’s spam experiences

What problems do ISPs have?

↳ Insecure customers

– very few real spammers sending directly !

- Botnets

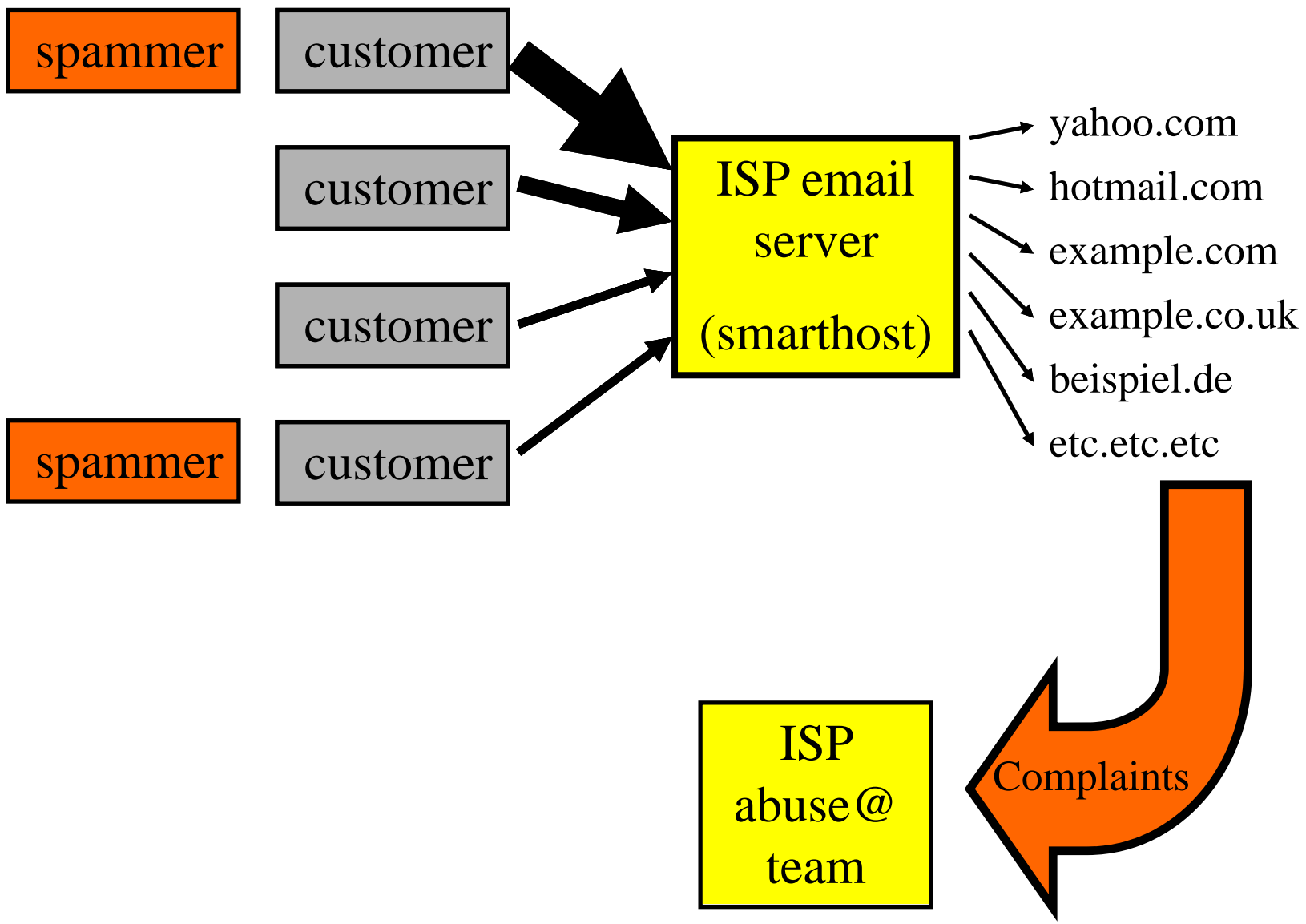
– compromised end-user machines

- SOCKS proxies &c

– mis-configuration

- SMTP AUTH

– Exchange “admin” accounts + *many others*

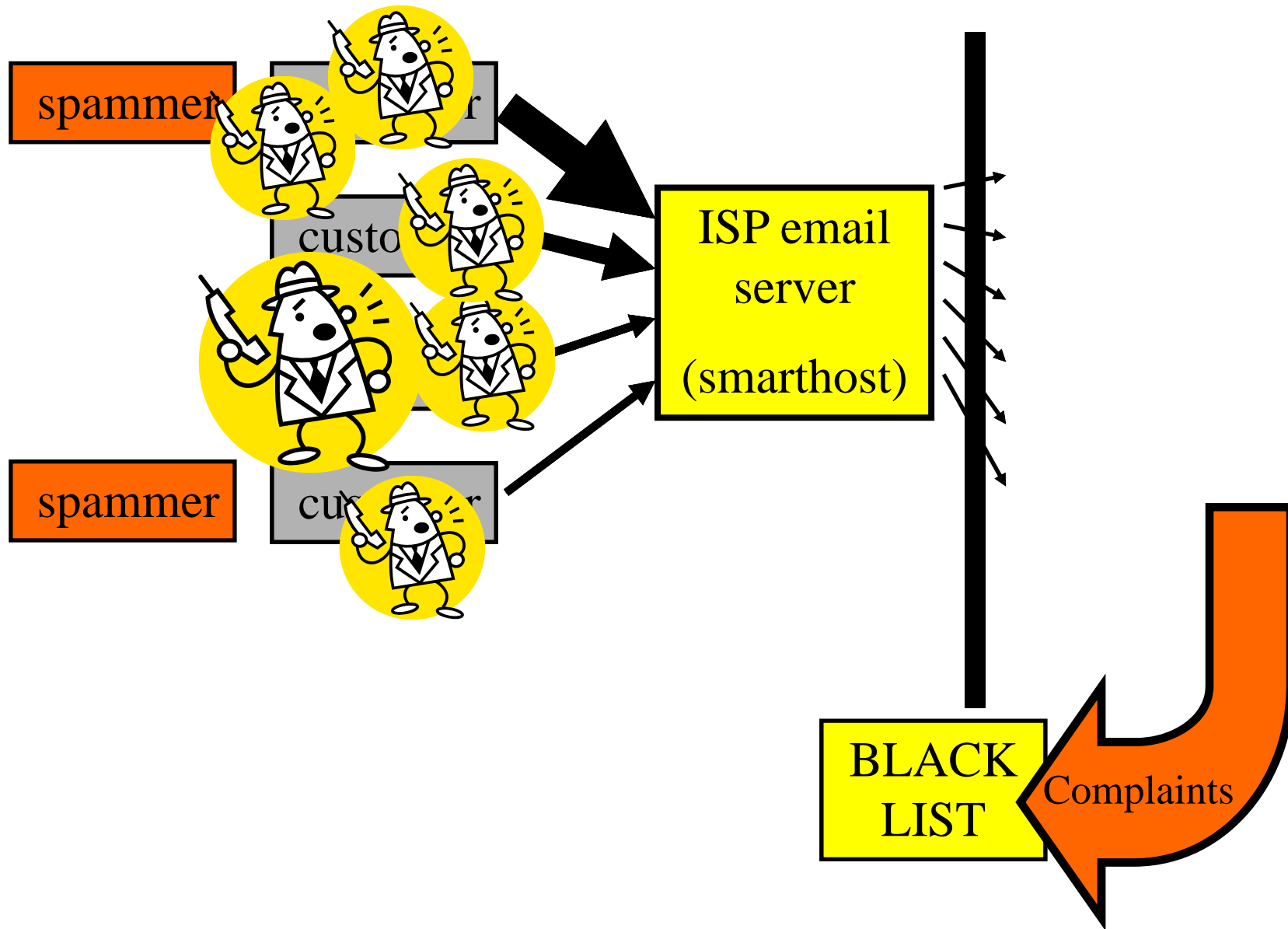


ISP's Real Problem

- Blacklisting of IP ranges & smarthosts
- Blocking by large email systems

HENCE:

- Rapid action necessary to ensure continued service to all other customers
- But reports may go to the blacklist and not to the ISP (or will lack essential details)

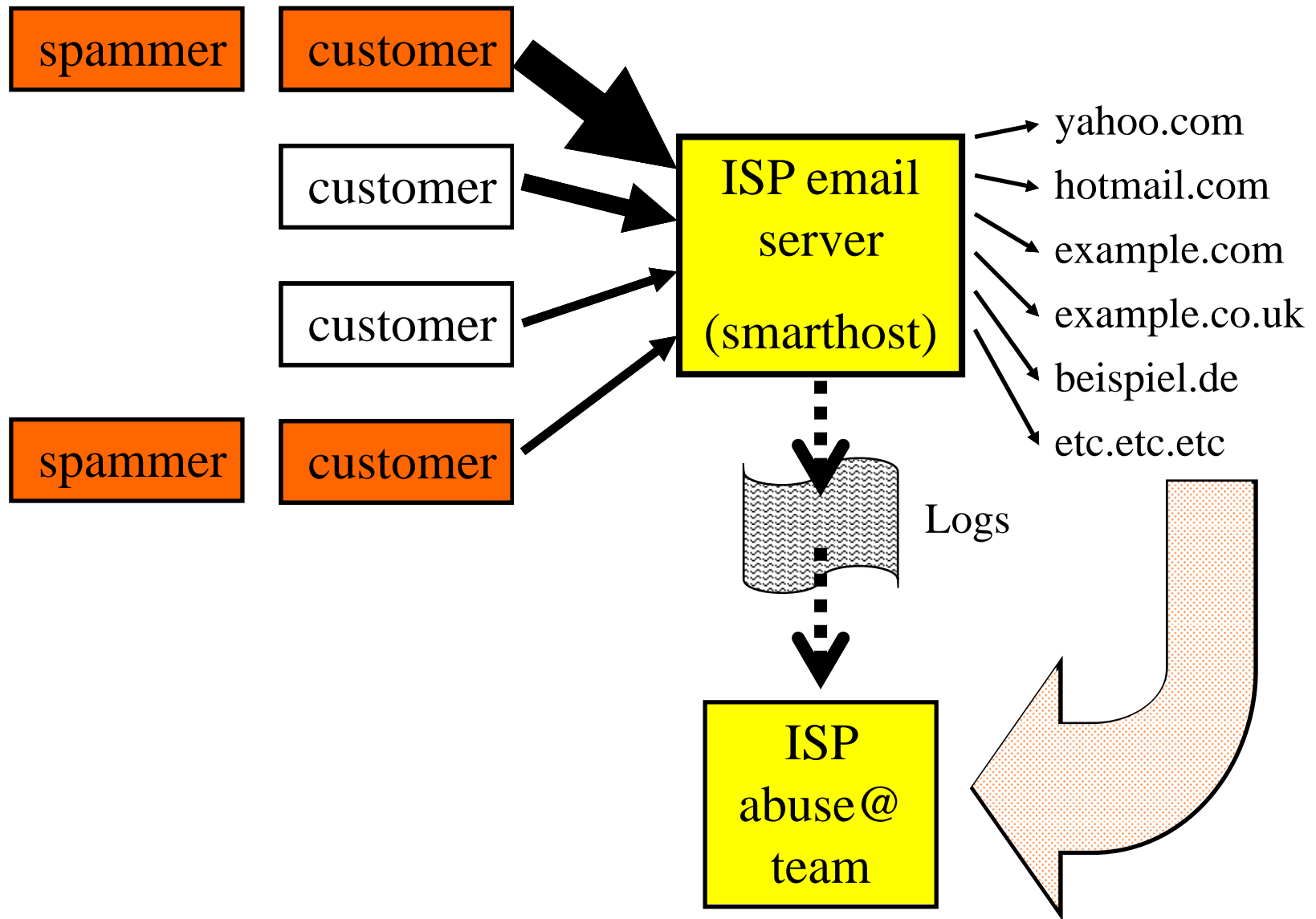


Spotting outgoing spam

- Expensive to examine outgoing content
- Legal/contractual issues with blocking
 - “false positives” could cost you customers
- Volume is not a good indicator of spam
 - many customers with occasional mailshots
 - daily limits only suitable for consumers
- “Incorrect” sender doesn’t indicate spam
 - many customers with multiple domains

Key insight (2003, still true)

- Lots of spam is to ancient email addresses
- Lots of spam is to invented addresses
- Lots of spam is blocked by remote filters (!)
- Can process server logs to pick out this information. Spam has many delivery failures whereas legitimate email mainly works



Log processing heuristics

- ↳ **Report “too many” failures to deliver**
 - more than 20 works pretty well
- Ignore “bounces” !
 - have null “< >” return path, these often fail
 - detect rejection daemons without < > paths
- Ignore “mailing lists”
 - most destinations work, only some fail (10%)
 - more than one mailing list is a spam indicator!

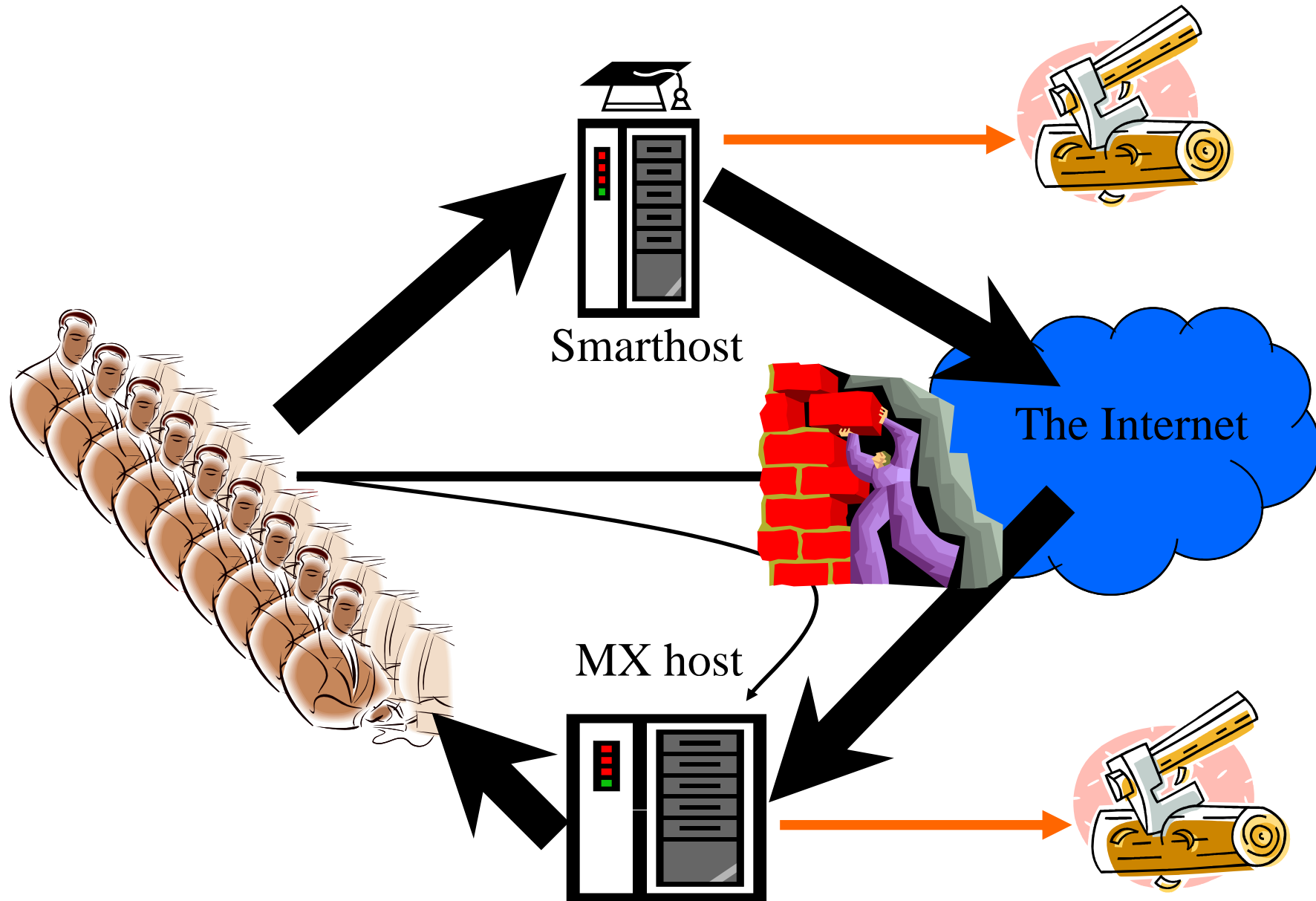
Bonus! also detects viruses

- Common for mass mailing “worms” to use address book (mainly valid addresses)
- But remote sites may reject malware
 - ALSO (and very useful) !
- Virus authors don’t know how to say HELO
 - or say HELO differently every time
- **So virus infections are also detected**
 - albeit, viruses less common these days

Bonus! can also detect loops

- Many people talk to themselves
 - e.g. unknown destinations sent to smarthost
- Many people's robots don't have null sender
 - vacation messages often have sender details
 - advert auto-responders want to be replied to
 - eventually these robots correspond with other dumbly configured systems and a mail loop is the result – sometimes of very high volume
- **Valuable to spot loops before 10K/day level!**

ISP email handling



Heuristics for incoming email

- Simple heuristics on failures work really well
 - just as for smarthost
- Multiple HELO lines very common
 - often match MAIL FROM (to mislead)
 - may match RCPT TO (? authenticator ?)
- Pay attention to spam filter results
 - but need to discount forwarding
- Outgoing email will fail on this machine

Spam being sent through the smarthost:

```
----- aafcu@office.com ->
2009-03-18 16:44:03 -> !aarond@unl.edu Size=1002
                also -> !aarond@unlserve.unl.edu
                -> aaronctidwell@yahoo.com
2009-03-18 16:44:06 -> aca@americancanoe.org Size=1000
                also -> aca@collegeofangiology.org
2009-03-18 16:44:11 -> acwriters@aol.com Size=1000
                also -> acwwa@hfx.andara.com
                -> aczesak@blaineds.org
2009-03-18 16:44:13 -> adrienne.shefik@dcsdk12.org Size=1000
                also -> adrianyearsley@yahoo.com
                -> adrielcg@respirnetpro.com
2009-03-18 16:44:24 -> afhe@primenet.com Size=1000
                also -> afhra.ahp@maxwell.af.mil
2009-03-18 16:44:25 -> !alamo_ccc@alamoccc.zzn.com Size=1000
                also -> !alamosa@fws.gov
                -> alameatoni@aol.com
2009-03-18 16:44:27 -> ags-registry@fao.org Size=1000
                also -> agstat@tds.net
                -> agthomson@msn.com
```

Outgoing email to the incoming email machine:

2009-03-10	18:38:39	muvt@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	20:05:41	tay@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	20:37:57	jip@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	20:38:54	tgp@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	21:10:14	dapum@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	22:14:46	dwd@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	22:47:01	xflj@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	22:47:58	llf@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	23:19:24	tnsk@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-11	23:52:33	bemb@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-12	00:23:59	bixfh@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-12	00:24:56	rqjan@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-12	00:56:18	nxf@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-12	00:57:15	stmx@	->	MATTASSOC1@mail.ru	unrouteable
2009-03-12	01:28:35	hxs@	->	MATTASSOC1@mail.ru	unrouteable

Varying HELO strings:

HELO = YJLBWOIVBH

2009-02-23 17:10:37 repliedlsoq@shoppingsingapore.com

-> haywood@let-it-be-thus.com

Size=1691

-> haywoodd@let-it-be-thus.com

-> hbxmyd@let-it-be-thus.com

-> healyn@let-it-be-thus.com

-> heardh@let-it-be-thus.com

-> heha@let-it-be-thus.com

HELO = FZNPWYWPF

2009-02-23 17:10:38 bridger@acetaxes.com

-> haven@let-it-be-thus.com

Size=1578

-> haynes@let-it-be-thus.com

-> haynesdd@let-it-be-thus.com

HELO = geos-ddce7df6b3

2009-02-23 19:45:46 emf_oohne@evenmorefun.com

-> d.levoi@evenoak.co.uk

Size=3520

Summary

- Processing outgoing server logs **works well**
 - keeps smarthosts out of blacklists
- Processing incoming server logs **effective**
 - little “looped back” traffic, but high signal to noise
- Production systems deployed at Demon Internet since September 2003, and continue in 2009 to be a major contributor to abuse reports
 - that’s a Good Thing!

<http://www.lightbluetouchpaper.org>

CEAS papers: <http://www.ceas.cc>

2004: Stopping spam by extrusion detection

2005: Examining incoming server logs

2006: Early results from spamHINTS

2007: Email traffic: A qualitative snapshot

2008: Do Zebras get more spam than Aardvarks?



**UNIVERSITY OF
CAMBRIDGE**
Computer Laboratory



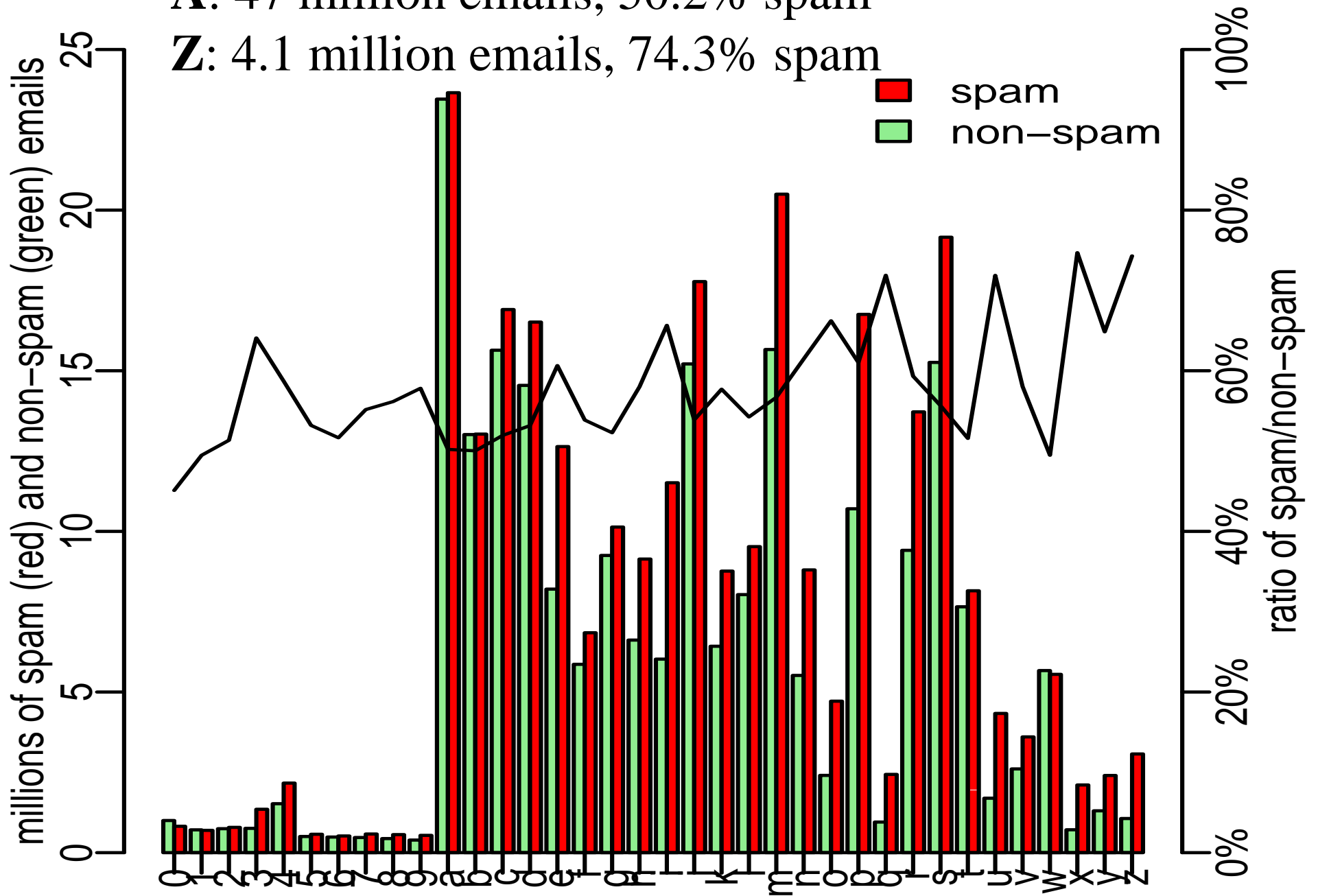
Demon

Demon email (Feb/Mar 2008)

- Ignored “bounces” (null sender)
 - mainly customer names taken in vain
- Treated n -addressed email as n emails
- 550 596 270 emails (8 million a day)
 - 56% were deemed to be spam by Cloudmark
- examined the first letter of the local parts
 - viz: was it addressed to an aardvark or a zebra

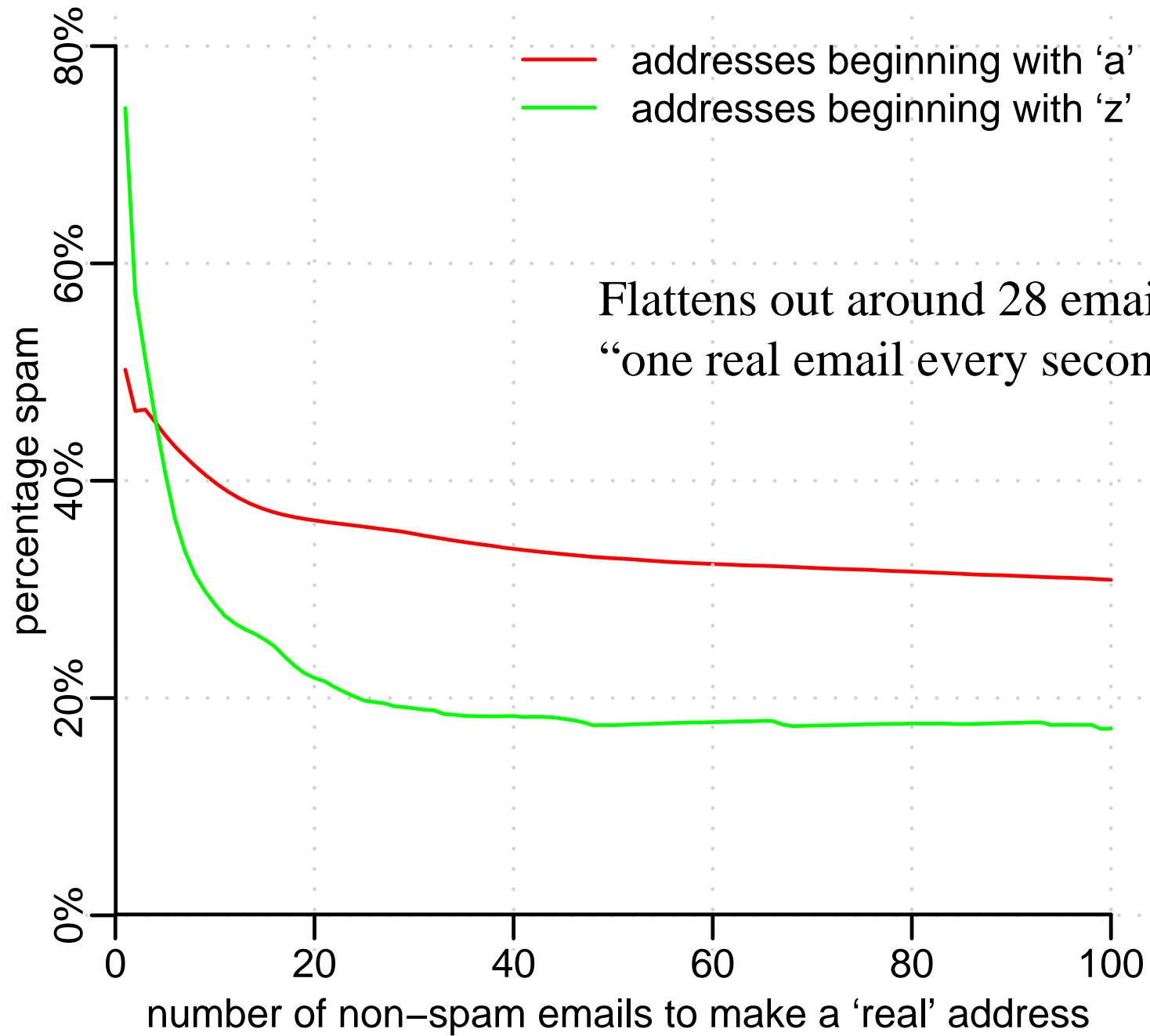
A: 47 million emails, 50.2% spam

Z: 4.1 million emails, 74.3% spam



“Real” Aardvarks/Zebras

- Not all email local parts are “real”
 - Demon doesn’t know a “ground truth”
 - non-real arise from “Rumpelstiltskin” or “dictionary” attacks... likely to be the underlying mechanism: your local part is guessed more often if there are a greater number of identical local parts
- So examine dataset to see which local parts receive n non-spam emails during the eight week period and deem these to be “real”



Results

- Zebras get way more spam than aardvarks
 - zebras 75%, aardvarks 50%
- But suppose we ignore imaginary animals
 - “real” zebras get 20% spam
 - whereas “real” aardvarks get 35% spam
- Filter designers might like to think about this
- Animals might like to consider a species change
- People might consider a new email address

<http://www.lightbluetouchpaper.org>

CEAS papers: <http://www.ceas.cc>

2004: Stopping spam by extrusion detection

2005: Examining incoming server logs

2006: Early results from spamHINTS

2007: Email traffic: A qualitative snapshot

2008: Do Zebras get more spam than Aardvarks?



**UNIVERSITY OF
CAMBRIDGE**
Computer Laboratory



Demon