# Do Zebras get more Spam than Aardvarks ?
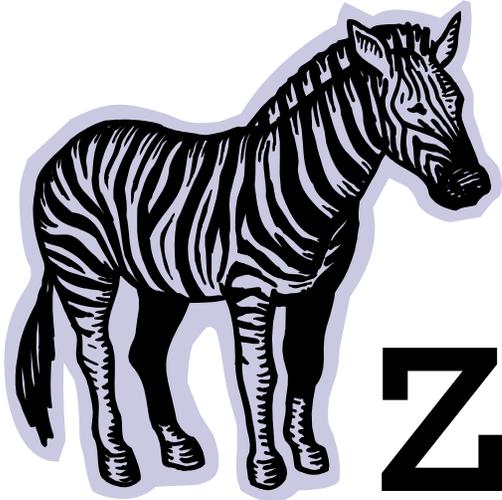
**Richard Clayton**

CEAS, Mountain View

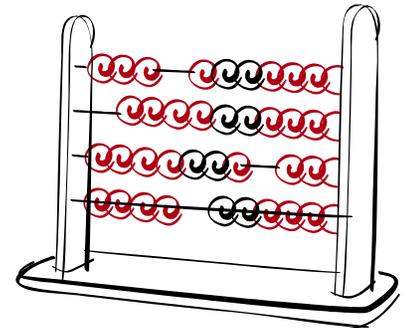22nd August 2008

UNIVERSITY OF
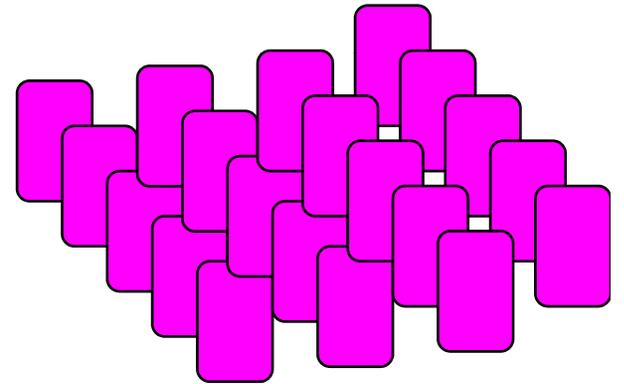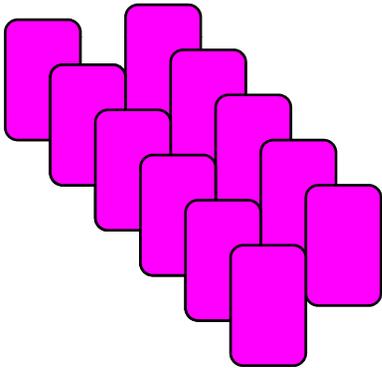CAMBRIDGE
Computer Laboratory

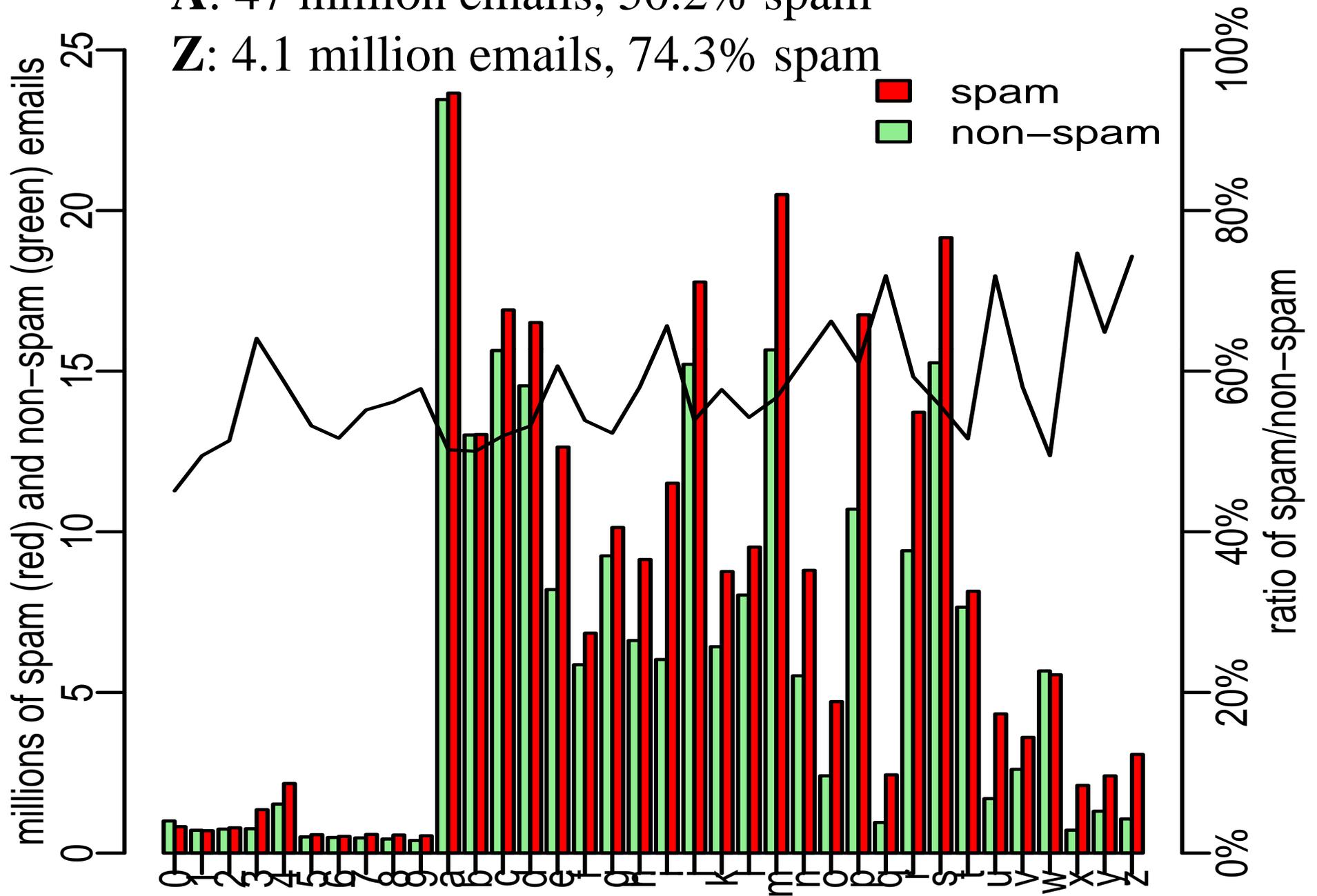Demon

Z

A

# Dataset

- Incoming email to Demon Internet
  - medium sized, long established UK ISP
  - c 150,000 customers, mainly ADSL, some dialup
  - mix of consumers, small & medium business
- Eight week dataset (1 Feb – 27 March 2008)
  - two public holidays (Easter)
  - cf CEAS 2007 which measured forwarding etc
  - BUT changes (PBL applied, ZEN greylisted)

# Raw numbers

- 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
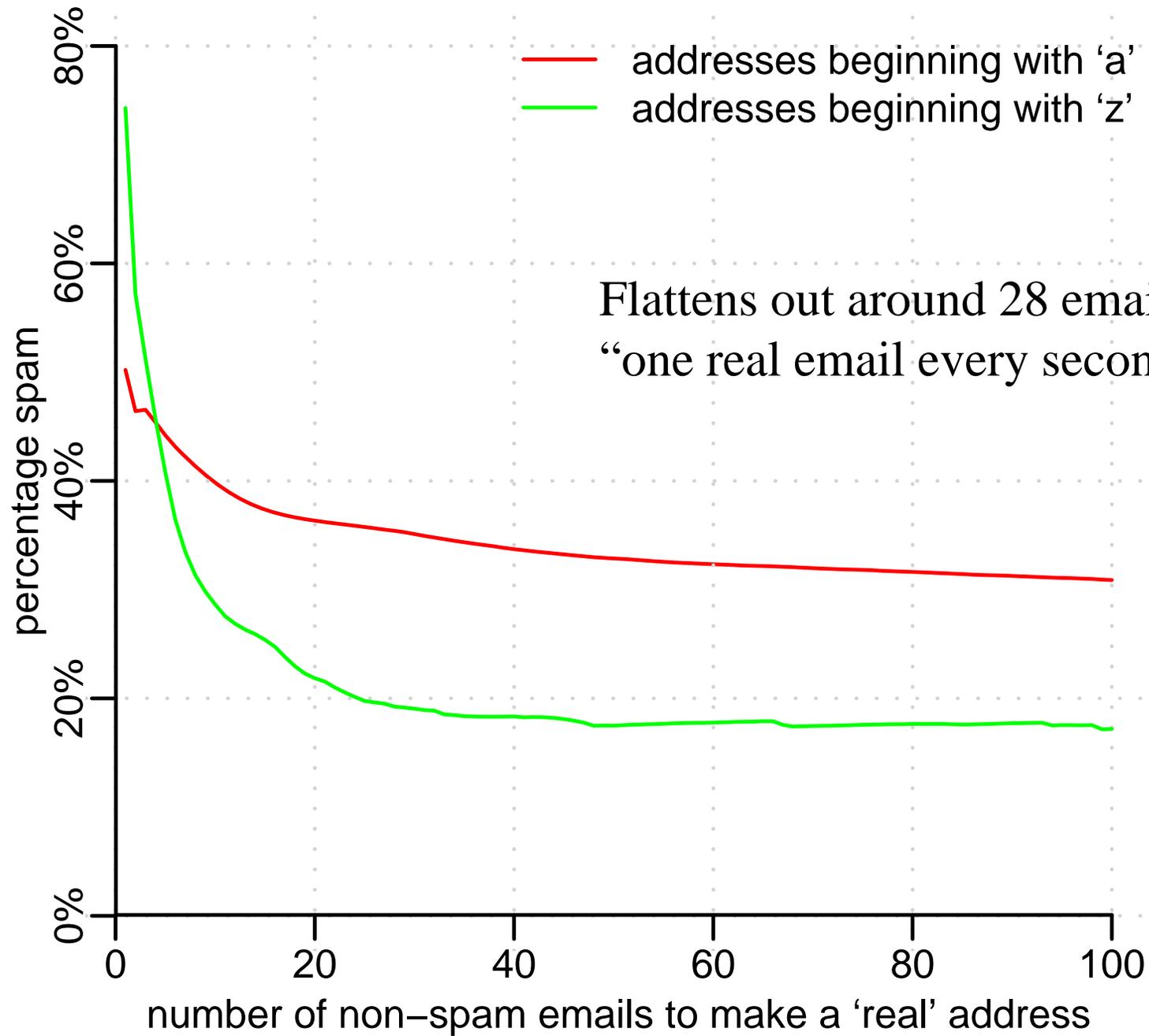  - viz: was it addressed to an <u>a</u>ardvark or a <u>z</u>ebra

**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 $n$ emails during the eight week period and deem these to be "real"
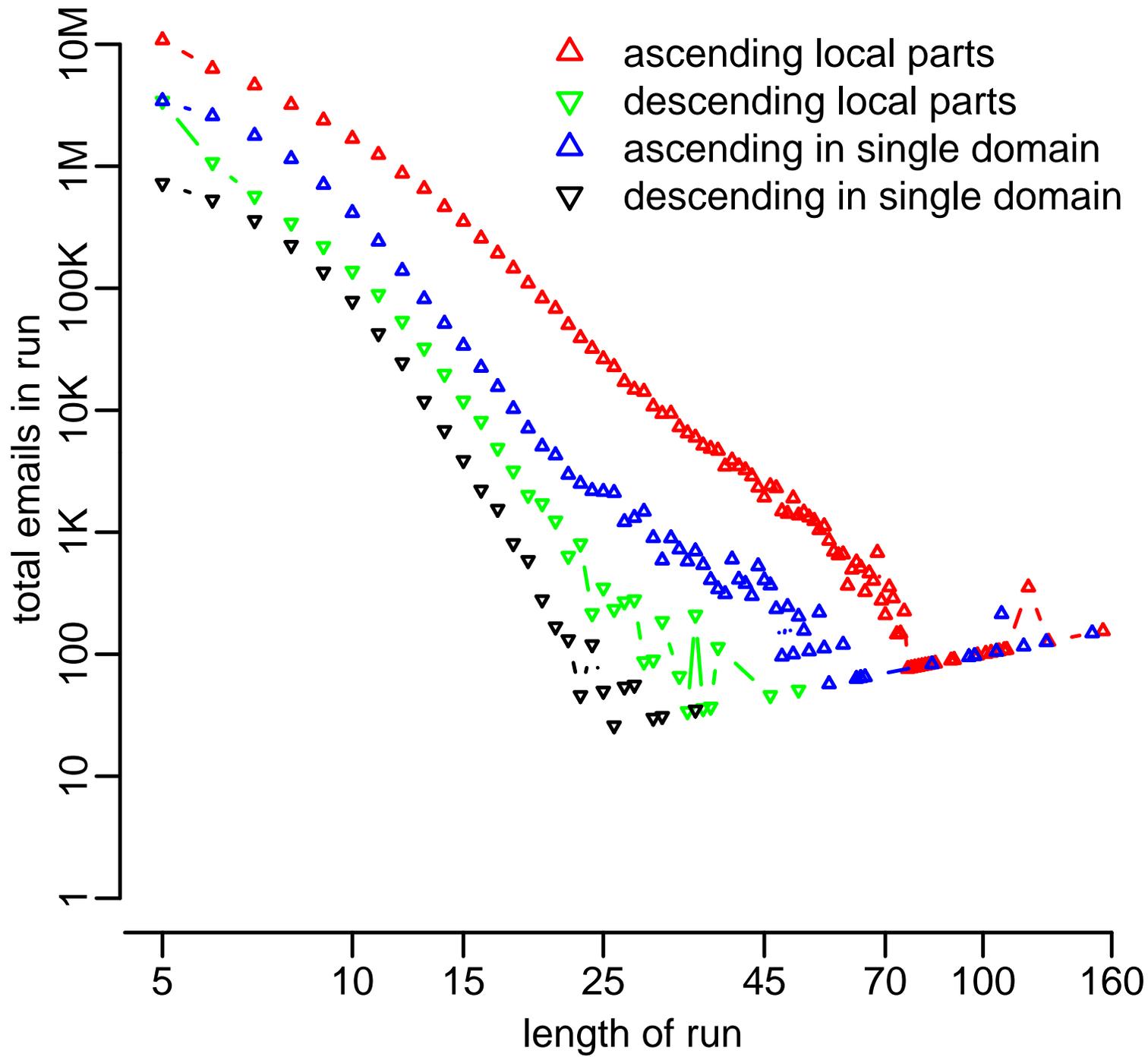
# Other amusement

- Can plot ratio of spam/ham for different starting letters
  - for example, "3" is a spam attractor
- Can use different definitions of what is "real" (for example 500+ non-spam emails)
  - see the paper (mercifully short!)

# Can we detect dictionary attacks?

- Expect to see "runs" of local parts in alpha order (ascending/descending)
- Might see "runs" across domains as well as within a single domain
- Evidence for these is unexpectedly weak:
  - Some runs of 100 or more
  - Only 2.9% of incoming spam in run of 5+

# Conclusions

- 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

# Do Zebras get more Spam than Aardvarks ?

**http://www.lightbluetouchpaper.org/**