

Alexander Serebrenik

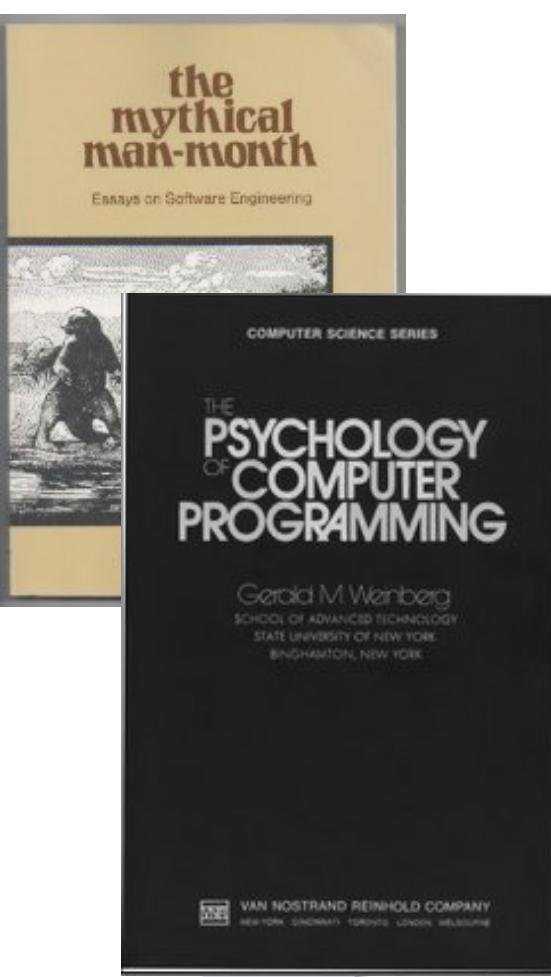
Eindhoven University of Technology, NL



Human Aspects of Software Engineering



And you!



1970s



open source

1998



MSR

SocialCom,
SocInfo, SBP,
CHASE

2008



2014

books,
conferences

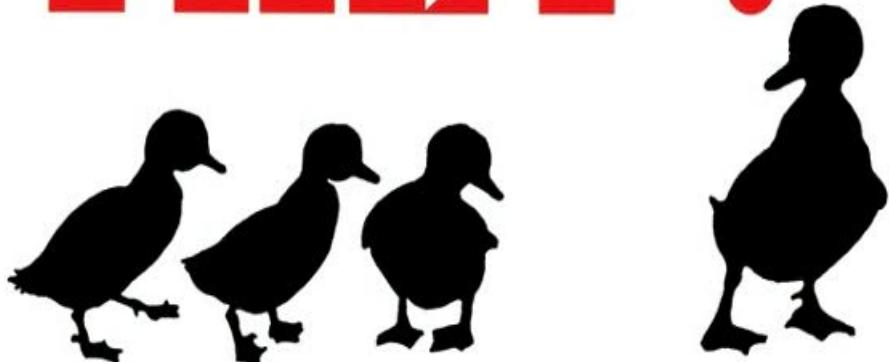


MS Windows post-release fault prediction	Precision predicted & correct / predicted	Recall predicted & correct / correct
Code churn	78.6%	79.9%
Code complexity	79.3%	66.0%
Code coverage	83.8%	54.5%
Code dependencies	74.4%	69.9%
Organizational structure	86.2%	84.0%
Socio-technical network	76.9%	70.5%





**WHO
ARE
THEY?**



TANA HOIBAN

**what
they do
in the
dark**

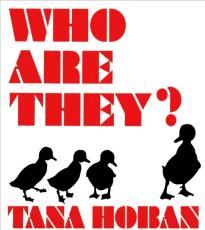
AMANDA COE



MEN



- > 90% in WordPress & Drupal
- > 95% in FLOSS surveys
- > 87% in GNOME
- > 70% in software-related jobs (NSF)



median



FLOSS
2013



3 years old

5-7

14-16

26-35

46-57

58-68

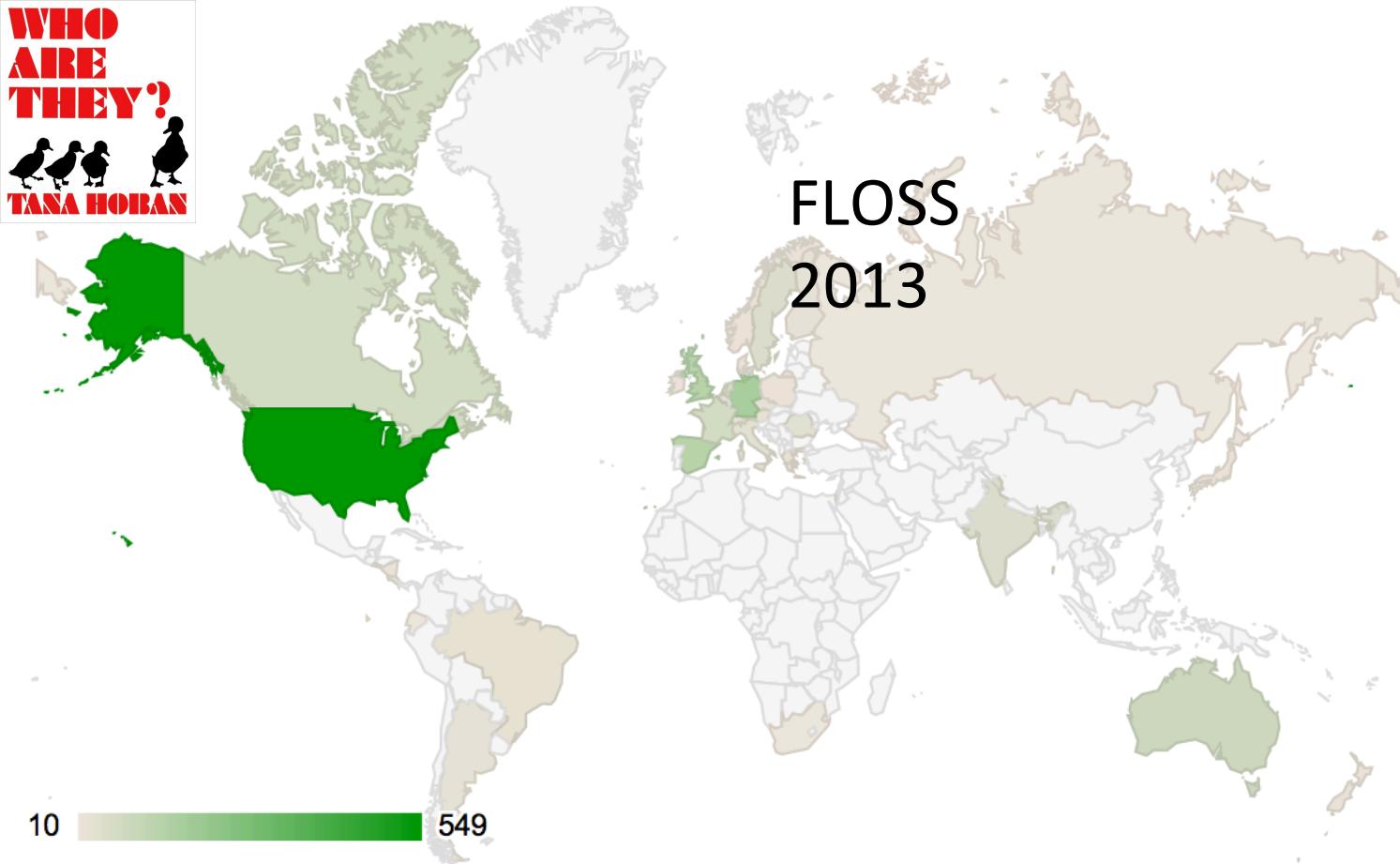
81-100

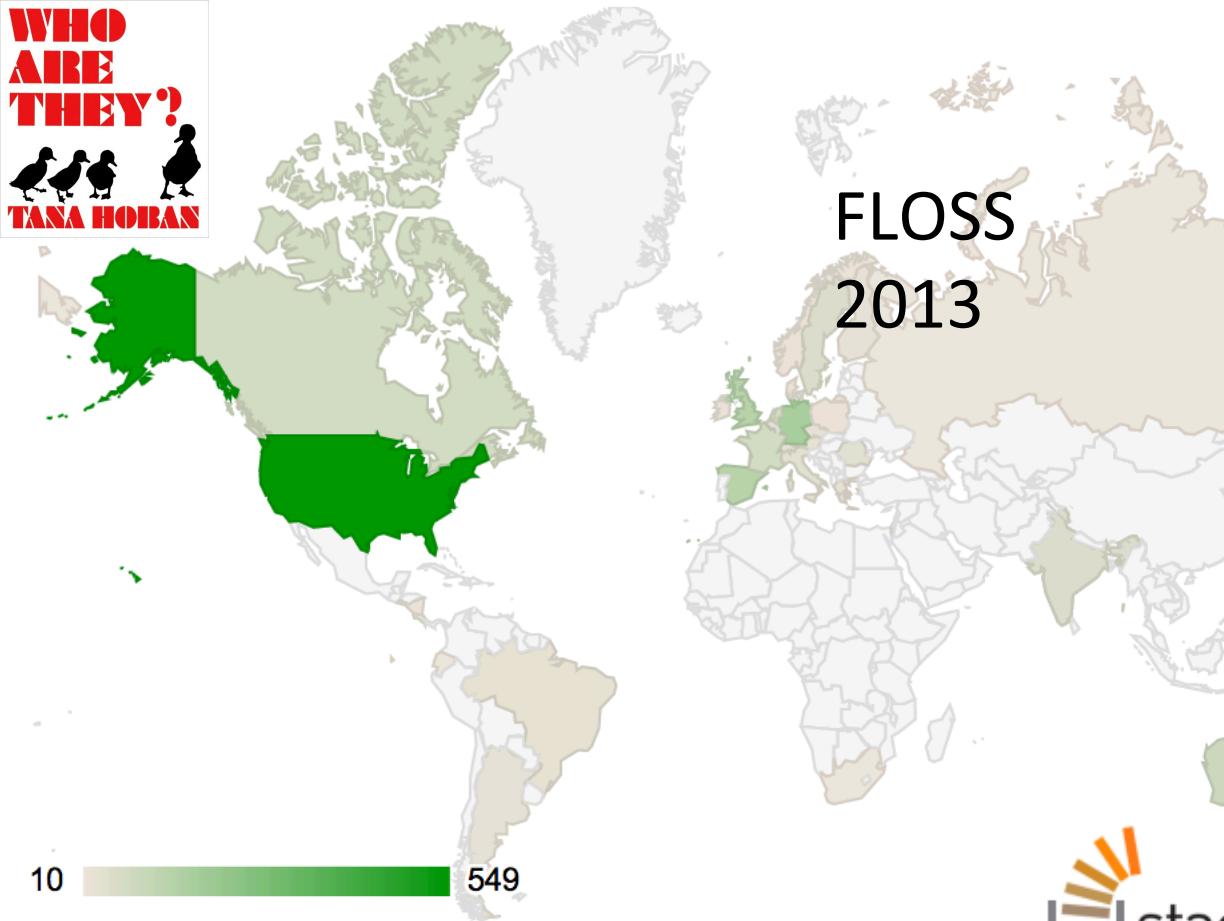
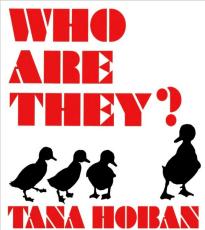
YOUNG(?)

Is Programming Knowledge Related to Age? An Exploration of Stack Overflow. Morrison, P., Murphy-Hill, E. MSR 2013. FLOSS 2013: A survey dataset about free software contributors: challenges for curating, sharing and combining. Robles, G., Arjona-Reina, L., Vasilescu, B., Serebrenik, A., Gonzalez-Barahona, J.M. MSR 2014

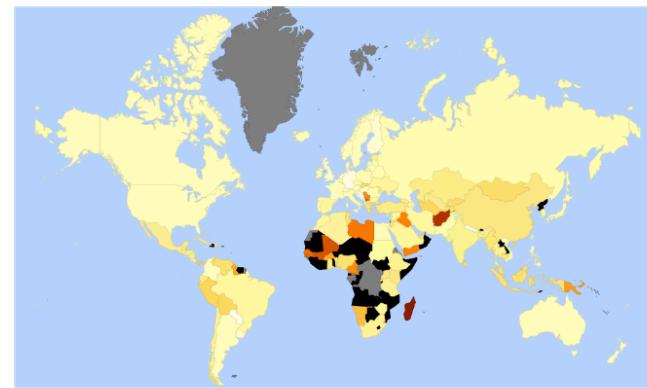
WHO
ARE
THEY?

TANA HOBAN

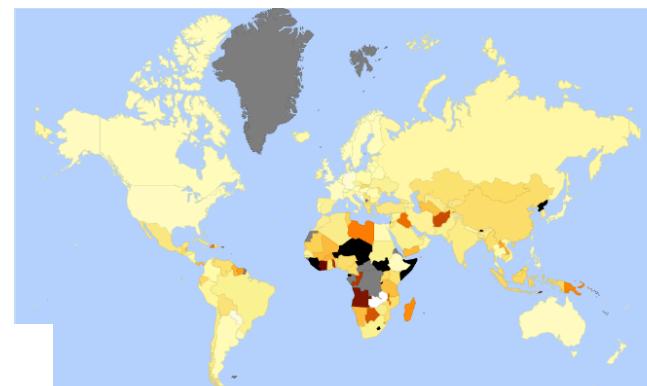




FLOSS
2013



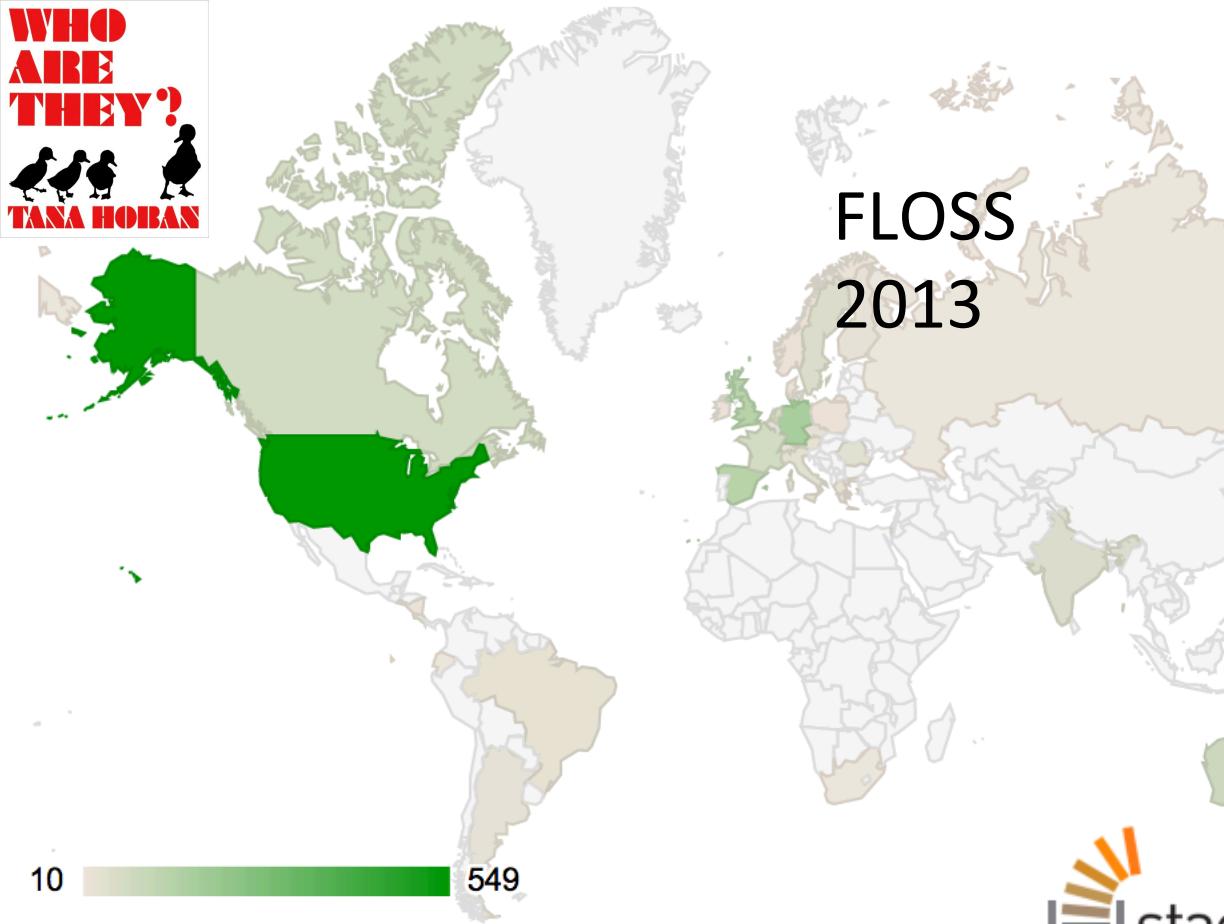
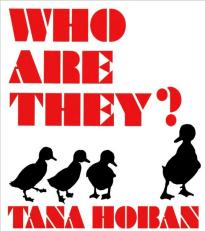
(d) Until June 2010



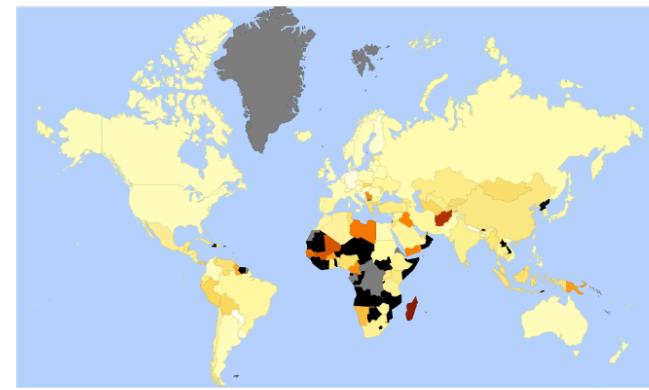
(g) Until December 2011



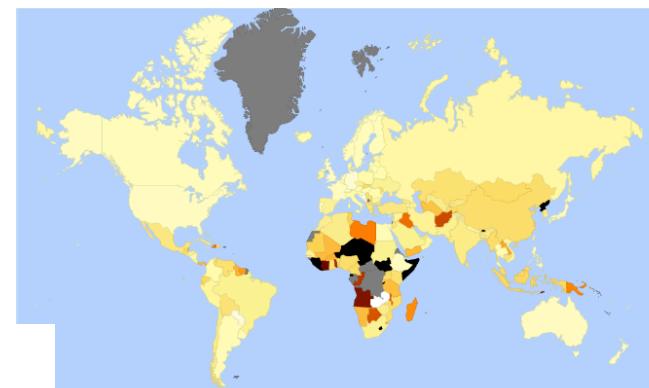
(j) Until June 2013



FLOSS
2013



(d) Until June 2010



(g) Until December 2011



Europe, US, CA, AU
Brazil/Argentina



(j) Until June 2013

what
they do
in the
dark



AMANDA COE





.cpp



.po



/test/



.jpg



/library/



.doc



makefile

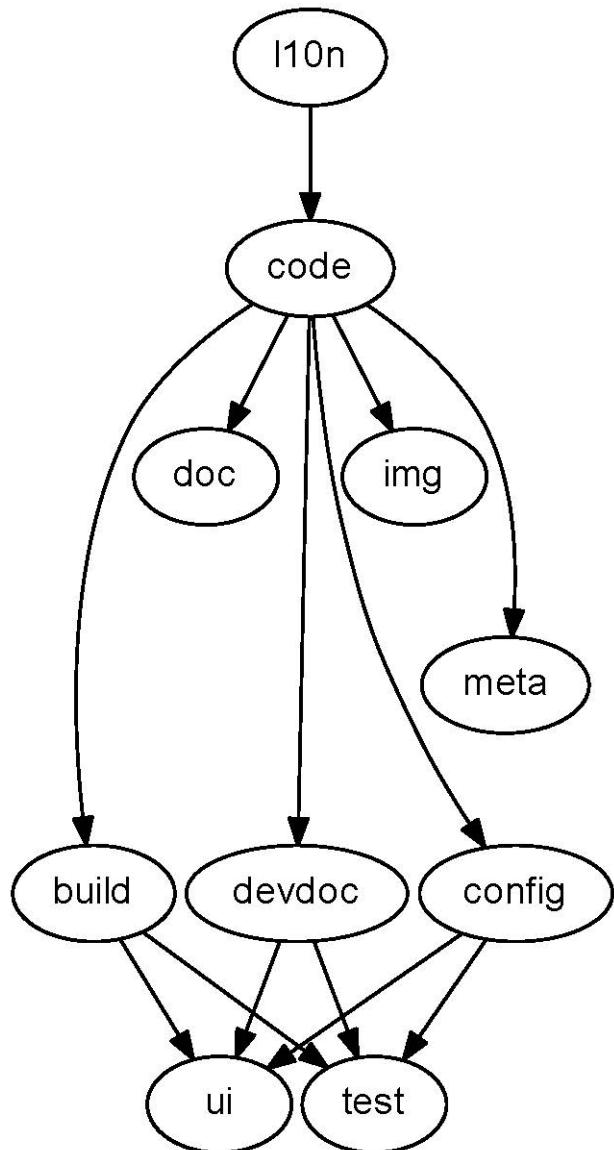


.sql

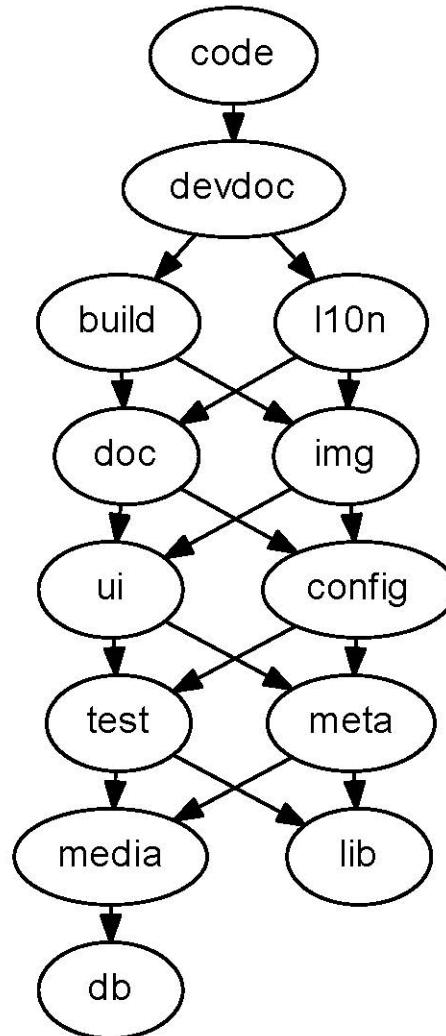


.conf

Occasional contributors



Frequent contributors



On the variation and specialization of workload---A case study of the Gnome ecosystem community Vasilescu, B., Serebrenik, A., Goeminne, M., Mens, T. Empirical Sw Engg

For which of those activities would
Joe Average be more active than
Jane Average?

For which activities no differences would
be observed?



.cpp



.po



/test/



.jpg



/library/



.doc



makefile



.sql



.conf



.cpp



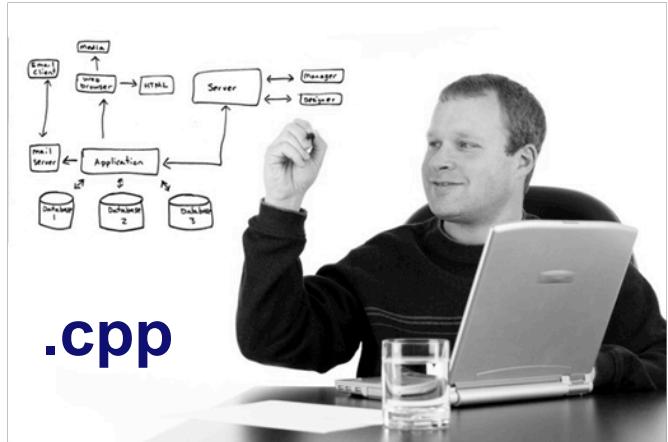
.po



.makefile

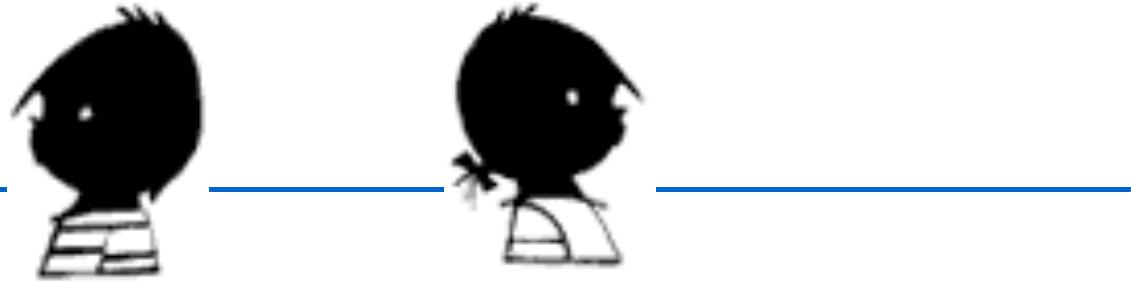


.sql



No differences
in *preferences*
for activities,
differences in
the *amount* of
commits





sample



WORDPRESS

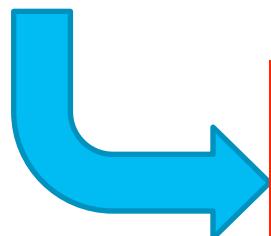


Drupal™

Engage
for longer

Ask more
questions

No diff in #answers



Women can
contribute to SO
but choose not to!



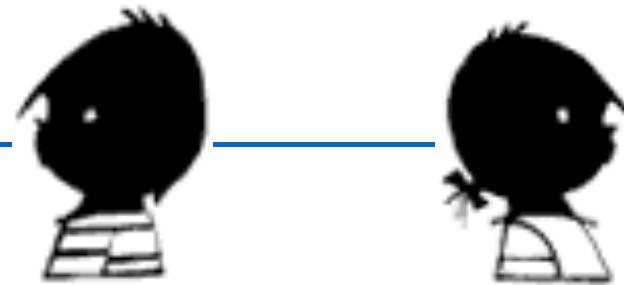
sample



WORDPRESS



Drupal™



No significant
differences in
#questions, #answers,
length of engagement



Disengagement of
women is activity/
platform specific!





- More active committers ask more on SO
- More active askers commit more on GitHub

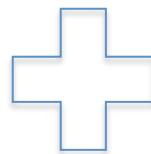


TABLE II: Mutual influence of StackOverflow activities and GitHub committing, for different committers (from least active $Q1$, to most active $Q4$).

Q	Influence of			
	asking on	committing on	answering on	committing on
committing	asking	committing	answering	
$Q1$	none	none	none	none
$Q2$	none	inconclusive	inconclusive	inconclusive
$Q3$	accelerates	accelerates	accelerates	accelerates
$Q4$	accelerates	accelerates	accelerates	accelerates



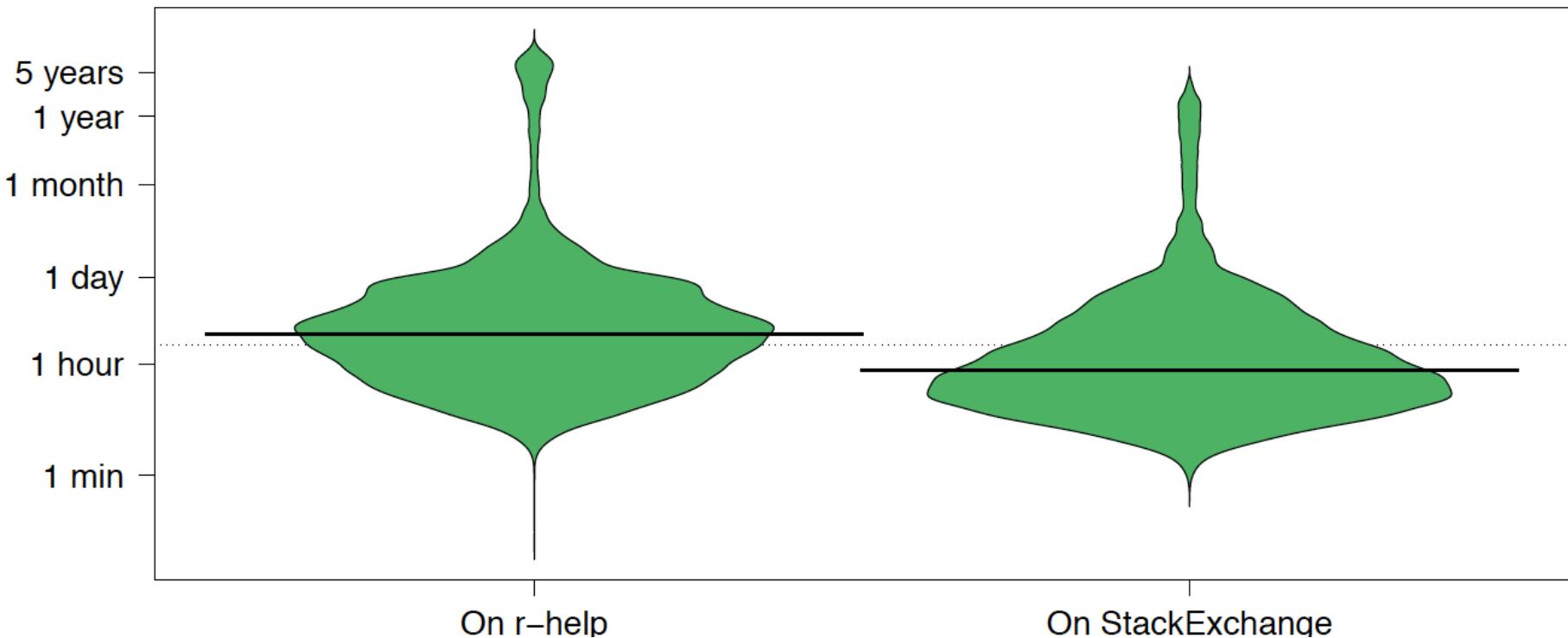
- Contributors to both communities (more likely developers than non-developers) are more active than those who focus on just one.



help



Speed of answers for r-help participants active on StackExchange





Follow on: [f](#) [t](#) [in](#) [+](#) [e](#)

Advertisement



Topics ▾

Reports ▾

Blogs ▾

Multimedia ▾

Tech Talk | At Work | Tech Careers

Older and Wiser... Up to a Point

By Philip E. Ross

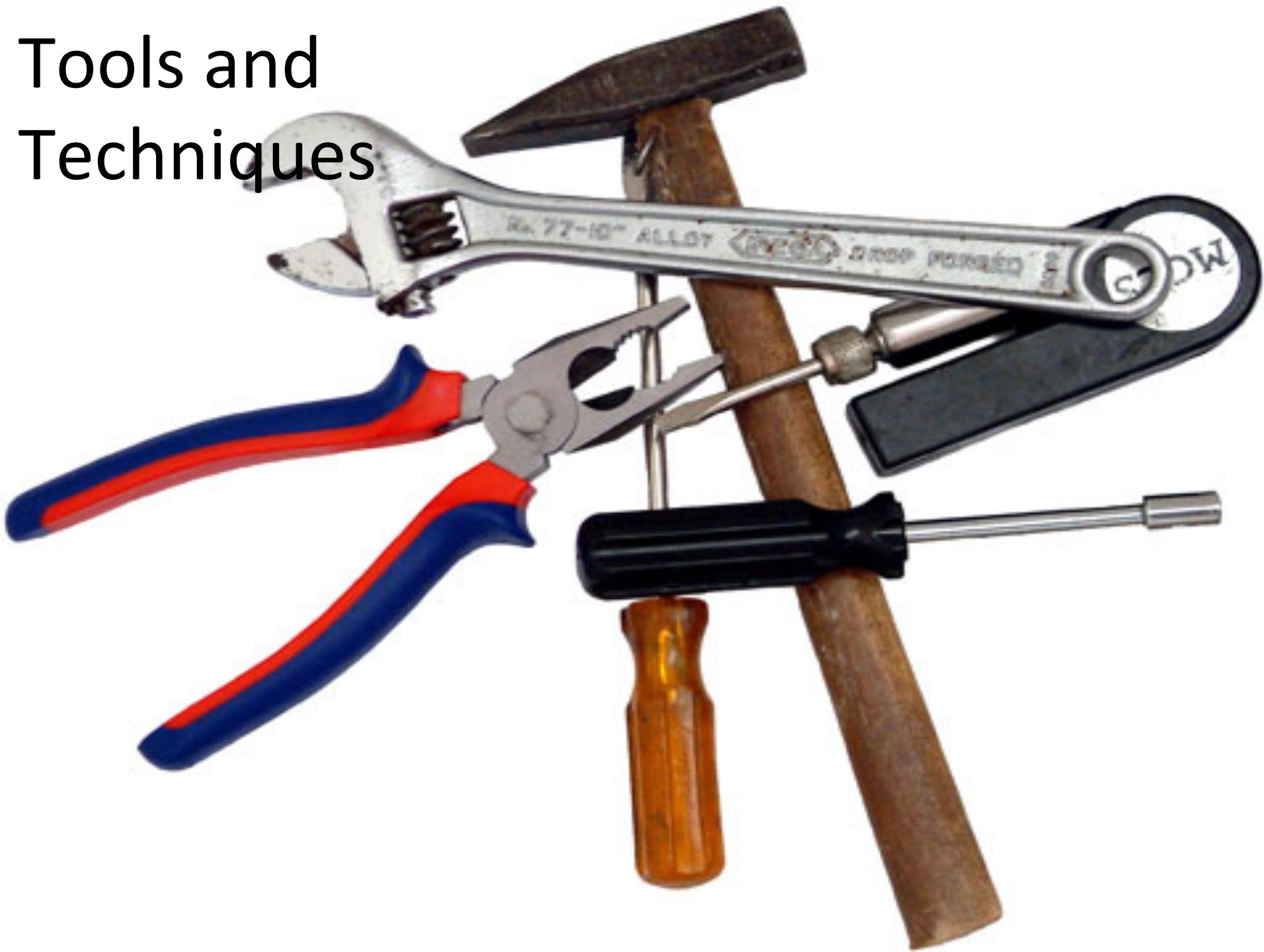
Posted 3 May 2013 | 11:21 GMT

Share |

Twitter

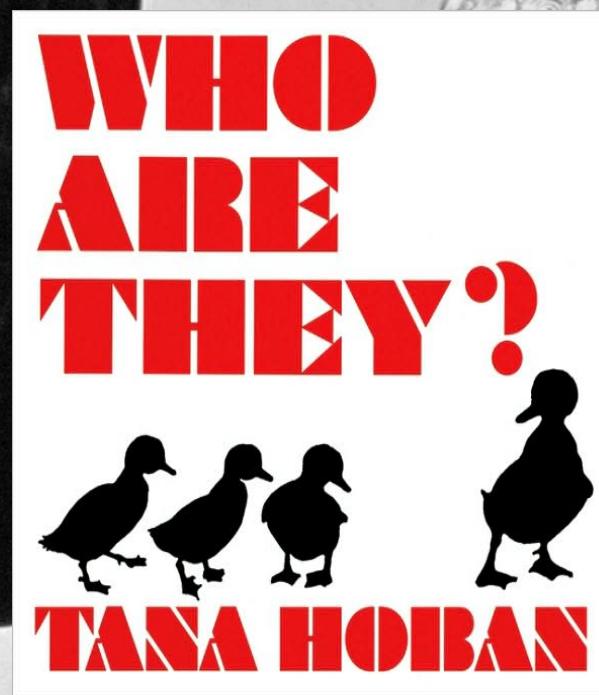
Facebook

Tools and Techniques



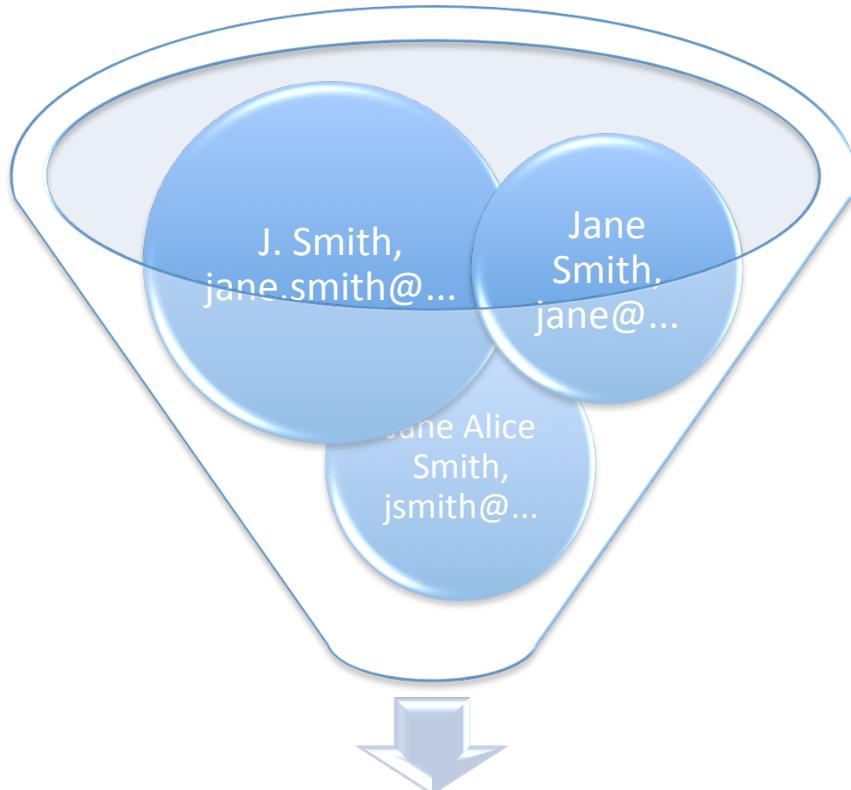


Robin Williams,
robinw@gmail.com



Euphegenia Doubtfire,
euphegenia@hotmail.com

Ordering	Rajesh Sola	Sola Rajesh
Spelling: misspelling, diacritics, punctuation	Rene Engelhard	Fene Engelhard
	Démurget	Demurget
	J. A. M. Carneiro	J A M Carneiro
Middle initials, patronyms, nicknames, additional surnames, incomplete names	Daniel M. Mueth	Daniel Mueth
	Alexander Alexandrov Shopov	Alexander Shopov
	Carlos Garnacho Parro	Carlos Garnacho
	Jacob “Ulysses” Berkman	Jacob Berkman
	A S Alam	Amanpreet Singh Alam
Name variants: transliteration, diminutives	Γιωργος	Georgios
	Mike Gratton	Michael Gratton
Software-specific: usernames, projects, tooling artefacts	mrhappypants	Aaron Brown
	Arturo Tena/libole2	Arturo Tena
	(16:06) Alex Roberts	Alex Roberts
Mix	Any combination of those	



Jane Smith

identity merging/
identity reconciliation/
developer matching

<Jane, jsmith@gmail.com>

<Jane Smith, jsmith@gmail.com>

- identical mail addresses → the same person
- also useful for MD5 hashes in Stack Overflow

<Jane Smith, jsmith@gmail.com>

<Jane Smith, jsmith@yahoo.com>

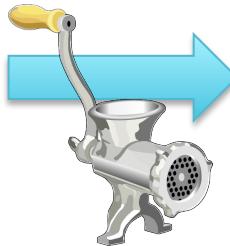
- likely to be the same person but might give false positives if these are different Janes...

Jane Smith

Latent Semantic Analysis

<John Doe,
<John Joseph Doe,

johnd@domainA>
johnd@domainA>



johnd@domainA:
{john, johnd,
joseph, doe}

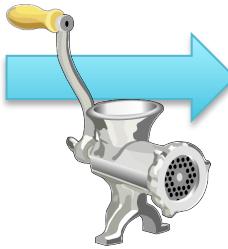
Document-term matrix

	johnd@...	j.doe@...	..
john	1
johnd	1
joseph	1
jdoe	?
doe	1

Latent Semantic Analysis

<John Doe,
<John Joseph Doe,

johnd@domainA>
johnd@domainA>



johnd@domainA:
{john, johnd,
joseph, doe}

Document-term matrix

	johnd@...	j.doe@...	..
john	1
johnd	1
joseph	1
jdoe	3/4
doe	1

max similarity(jdoe,
{john, johnd, joseph, doe})
= similarity(jdoe, doe)
= $1 - \text{Levenshtein}(jdoe, doe) / \max(\text{length}(jdoe), \text{length}(doe))$
= $1 - 1/4 = 3/4$

Latent Semantic Analysis

	johnd@..	j.doe@..	..
john	1
johnd	1
joseph	1
jdoe	3/4
doe	1

Inverse document frequency



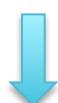
(Singular value decomposition)



Rank (noise) reduction

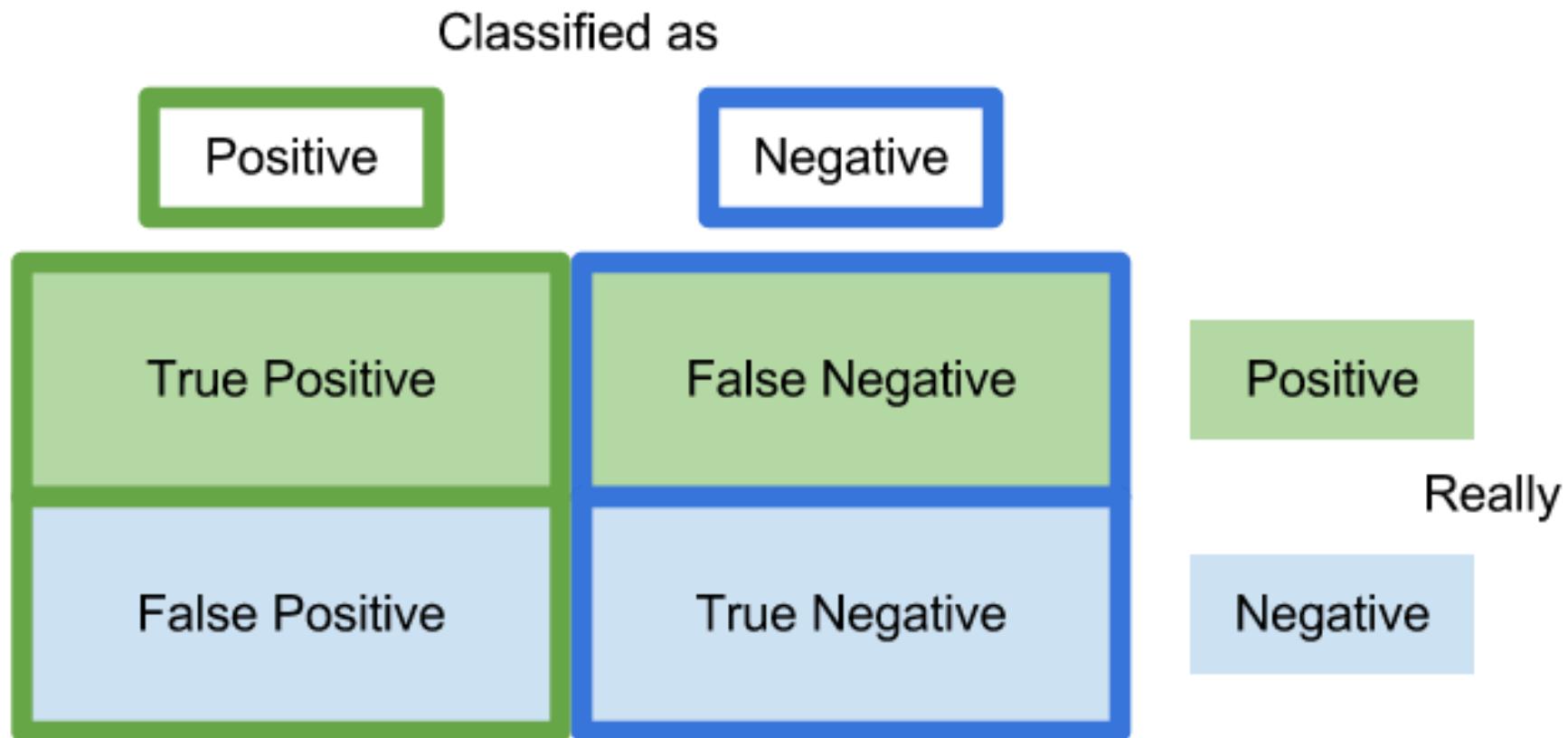


Cosine between documents

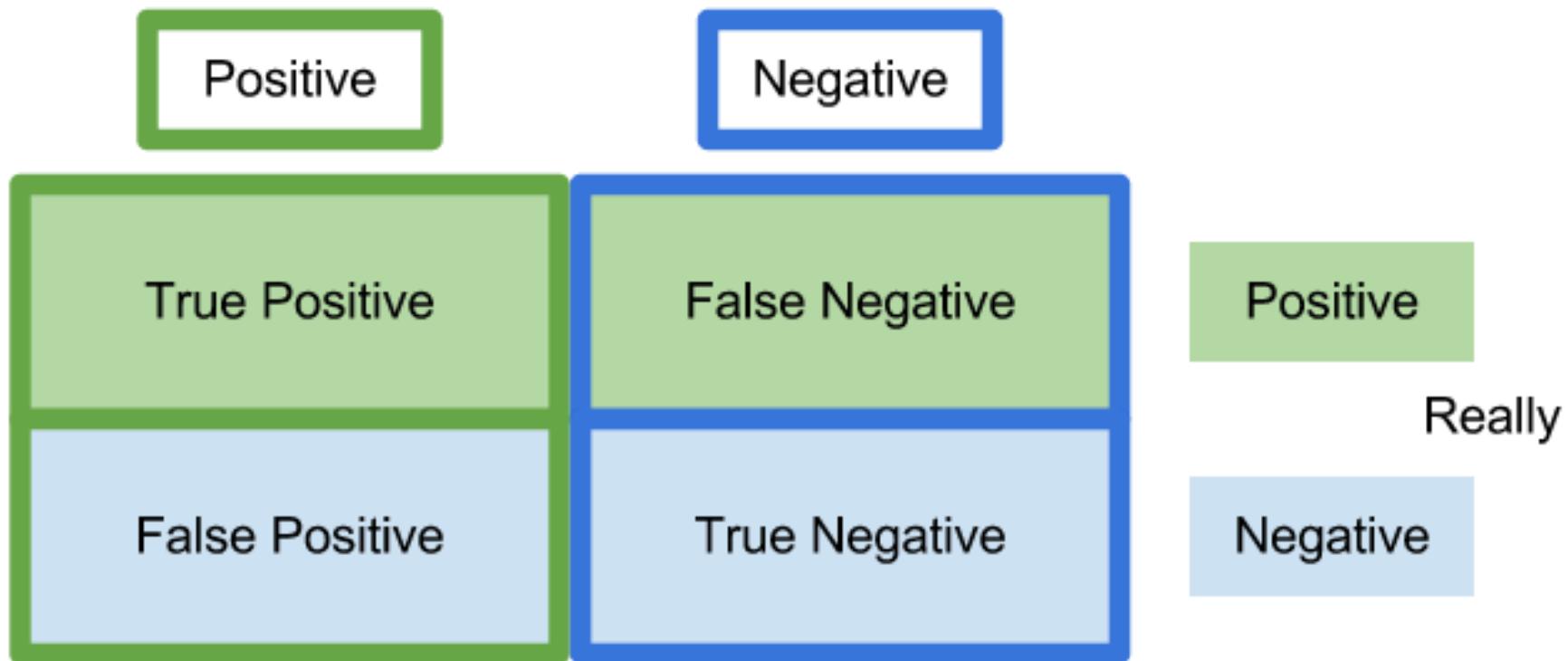


Merge similar documents

Which technique is better?



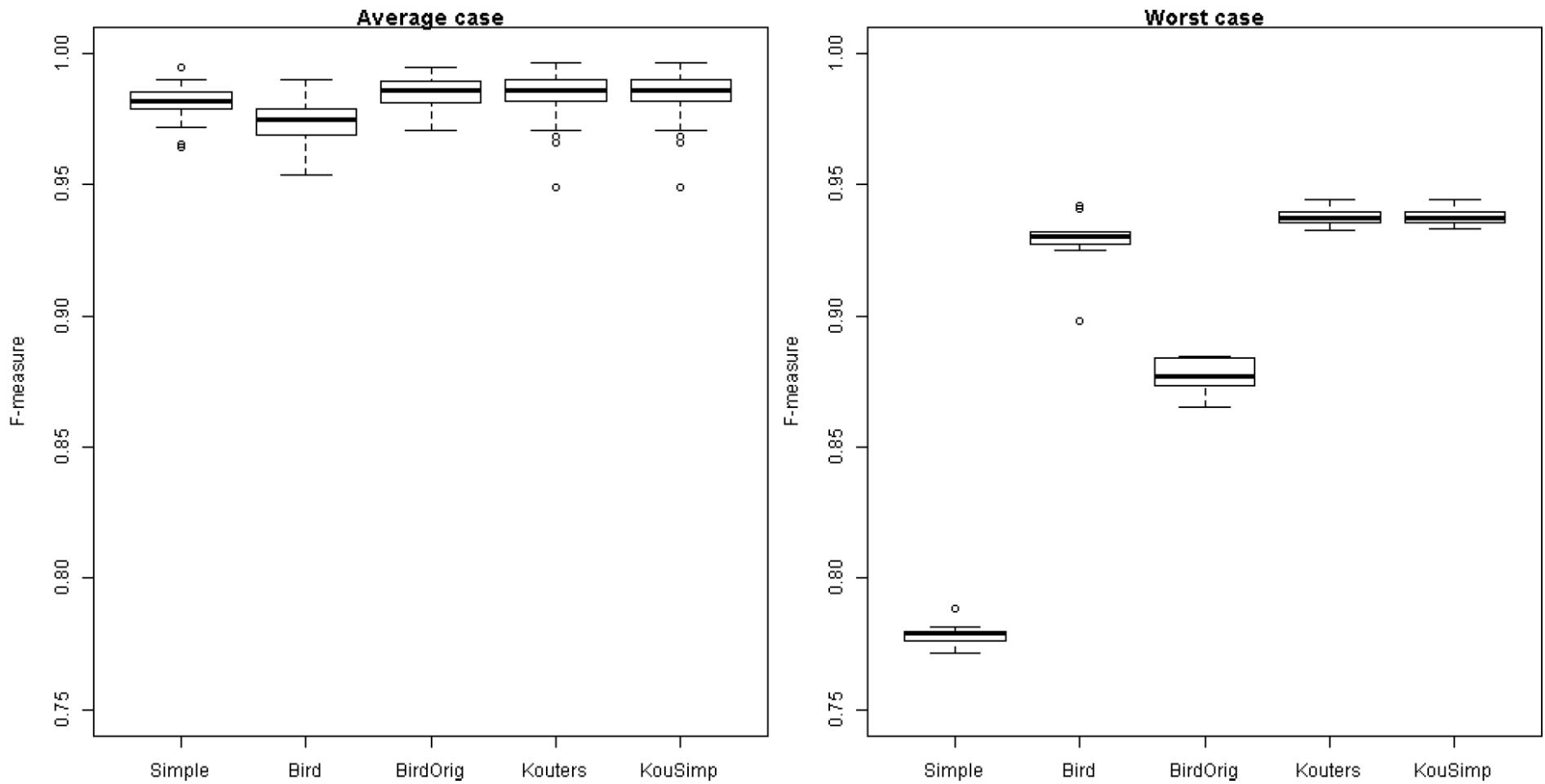
Classified as



$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$



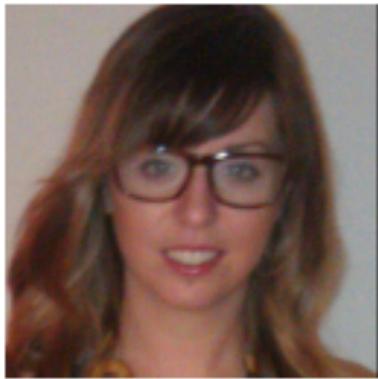
Most contributors have only one “identity”: all techniques
are good enough for the average case.
Advanced techniques are better in presence of noise.
Choose your weapons wisely!

A dark, moody background featuring a person's silhouette against a bright, cloudy sky.

What is
your
gender?

Sara Chipps

[less info](#)



[bio](#)

[visit:](#)

4,168

reputation

[stats](#)

● 5 ● 35 ● 83

Sara Chipps

[less info](#)



[bio](#)

[visit:](#)

4,168

reputation

• 5 • 35 • 83

Andrea Ambu

[less inf](#)



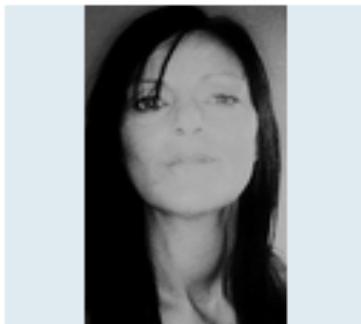
[bio](#)

[visi](#)

Andrea Smith

[sta](#)

[less](#)



[t](#)

[v](#)

[s](#)

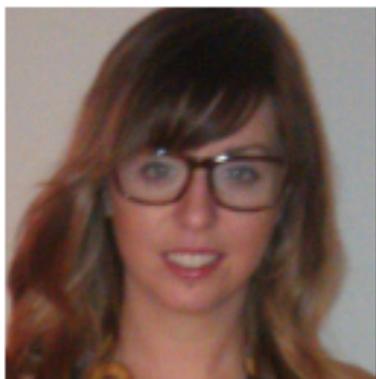
1

reputation

?

Sara Chipps

[less info](#)



[bio](#)

[visit](#)

4,168

reputation

• 5 • 35 • 83

Andrea Ambu

[less info](#)



[bio](#)

[website](#)

[andreaa.com](#)

[location](#)

[Italy](#)

[age](#)

[25](#)

[visits](#)

[member for](#)

[5 years, 1 months](#)

[seen](#)

[Oct 20 at 11:01](#)

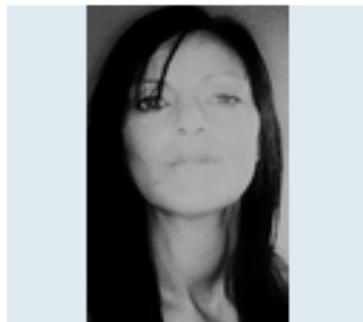
[state](#)

[less in](#)

[profile views](#)

[1,232](#)

Andrea Smith



[t](#)

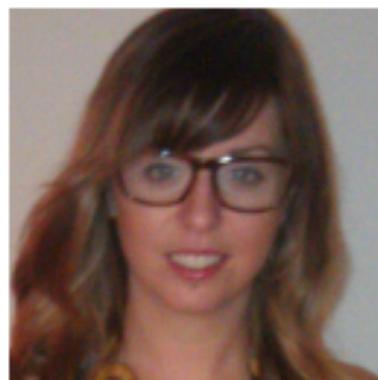
[v](#)

[s](#)

1

reputation

Sara Chipps less info



[bio](#)

[visit](#)

4,168
reputation

• 5 • 35 • 83

Andrea Ambu less info



[bio](#)

[website](#)

[andreaa.com](#)

[location](#)

[Italy](#)

[age](#)

[25](#)

[visits](#)

[member for](#)

[5 years, 1 months](#)

[seen](#)

[Oct 20 at 11:01](#)

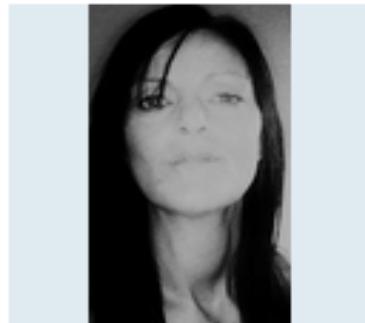
[state](#)

[less in](#)

[profile views](#)

[1,232](#)

Andrea Smith



[t](#)

[v](#)

[s](#)

1
reputation

Name +
Location =
Gender

vsushkov

[less info](#)



1,678
reputation

• 1 • 5 • 15

[bio](#)

[location](#)

[age](#)

[website](#)

[vsushkov.com](#)

[Taganrog, Russia](#)

[visits](#)

[member for](#)

3 years, 3 months

[seen](#)

15 hours ago

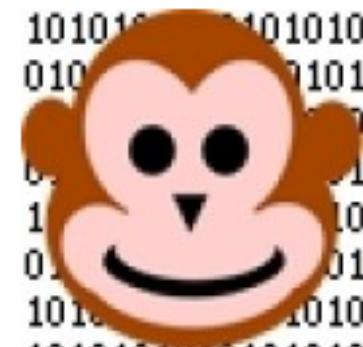
[stats](#)

[profile views](#)

188

w35l3y

[less info](#)



1,908
reputation

• 9 • 27

w35l3y ⇒ wesley

Lonzo

[less info](#)



[bio](#)

[visits](#)

[stats](#)

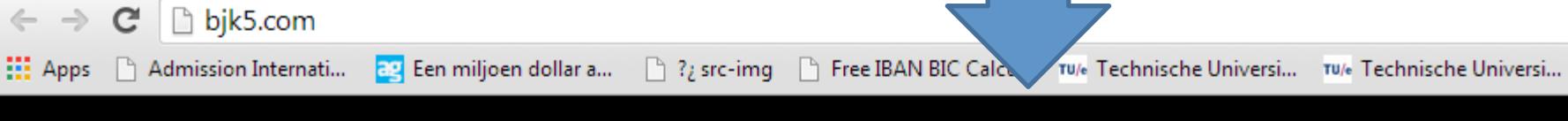
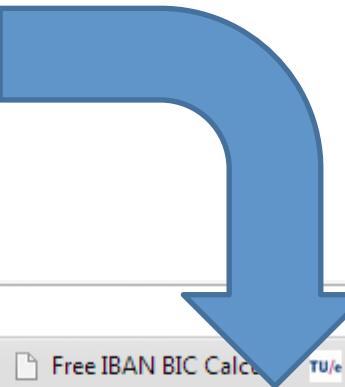
1,177
reputation

• 3 • 12 • 18

Name +
Location =
Gender



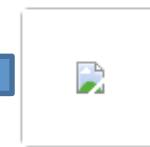
<i>bio</i>	website	bjk5.com
	location	United States
	age	30
<i>visits</i>	member for	5 years, 2 months
	seen	21 hours ago



Heuristics:

title + first h1

<title>Ben Kamens</title>
...
<h1>We're willing to be embarrassed about what we *haven't* done</h1>



Ben Kamens

is lead dev at [Khan Academy](#), and has been a proud part of [Fog Creek](#)

Ben Kamens We're willing to be embarrassed about what we haven't done...

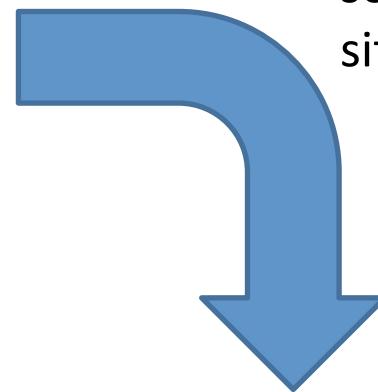
Stanford Named Entity Tagger

<PERSON>Ben Kamens</PERSON>
We're willing to be embarrassed about what we haven't done...

Quality of gender resolution: Survey

Self-identification	As inferred			Total
	M	F	?	
M	60	3	43	106
F	2	5	4	11

+ avatars, other
social media
sites (manually)

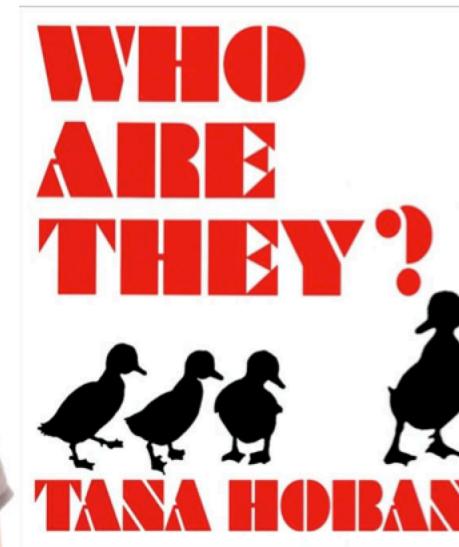


Self-identification	As inferred			Total
	M	F	?	
M	90	3	13	106
F	2	9	0	11

SUMMARY



Tools and Techniques



SANER 2015

22nd International Conference on Software Analysis, Evolution and Reengineering

- Abstracts: November 7, 2015
- Papers: November 14, 2015