## Who's who (in GNOME)?

e.t.m.kouters@student.tue.nl + b.n.vasilescu@tue.nl + a.serebrenik@tue.nl + m.g.j.v.d.brand@tue.nl Eindhoven University of Technology, The Netherlands

## Software ecosystems

- different projects
- different repositories per project (version control systems, mailing lists, ...)
- same community of developers

## Multiple aliases

Contributors sign off using a <name, email> alias. The same person often uses multiple aliases.

Names: John Travolta



**Emails**: j.travolta@domainA





- Travolta John
- J Travolta
- John
- John Trabolta
- John J. Travolta
- John Joseph Travolta
- John "Bone" Travolta
- john.travolta@domainB john DOT travolta AT
  - domainC
  - jtravolta@domainD
  - john@domainE

Identity merge algorithms

Merge aliases belonging to the same "real" person.

• Goeminne & Mens (2011) – simple, state of the art 〈John Travolta, j.travolta@domainA>

Bird et al. (2006) – more complex, inspiring

dynamical Travolta, bone@domainA>

dJohn Travolta, john@domainB>

j.travolta@domainA> 🕥 <John Travolta, <John J. Travolta, john@domainB>



<John Travolta, <John F. Kennedy, john@domainB>

john@domainA>



<John,







<John Travolta, john@domainA> <Travolta John, travolta@domainB> John !~Levenshtein Travolta



Latent Semantic Analysis

johnt@domainA> <John Travolta, <John Joseph Travolta, johnt@domainA>



johnt@domainA: {john, johnt, joseph, travolta}

## Results for GNOME

- 8618 different aliases
- only 4989 unique!

1.00





max similarity(jtravolta, {john, johnt, joseph, travolta}) = similarity(jtravolta, travolta) = 1 - Levenshtein(jtravolta, travolta) = 1 - 1/9 = 8/9





/ Department of Mathematics and Computer Science