

## 2 From open-source software to Wikipedia: ‘Backgrounding’ trust 3 by collective monitoring and reputation tracking

4 Paul B. de Laat

5  
6 © Springer Science+Business Media Dordrecht 2014

7 **Abstract** Open-content communities that focus on co-  
8 creation without requirements for entry have to face the  
9 issue of institutional trust in contributors. This research  
10 investigates the various ways in which these communities  
11 manage this issue. It is shown that communities of open-  
12 source software—continue to—rely mainly on hierarchy  
13 (reserving write-access for higher echelons), which sub-  
14 stitutes (the need for) trust. Encyclopedic communities,  
15 though, largely avoid this solution. In the particular case of  
16 Wikipedia, which is confronted with persistent vandalism,  
17 another arrangement has been pioneered instead. Trust (i.e.  
18 full write-access) is ‘backgrounded’ by means of a per-  
19 manent mobilization of Wikipedians to monitor incoming  
20 edits. Computational approaches have been developed for  
21 the purpose, yielding both sophisticated monitoring tools  
22 that are used by human patrollers, and bots that operate  
23 autonomously. Measures of reputation are also under  
24 investigation within Wikipedia; their incorporation in  
25 monitoring efforts, as an indicator of the trustworthiness of  
26 editors, is envisaged. These collective monitoring efforts  
27 are interpreted as focusing on avoiding possible damage  
28 being inflicted on Wikipedian spaces, thereby being  
29 allowed to keep the discretionary powers of editing intact  
30 for all users. Further, the essential differences between  
31 backgrounding and substituting trust are elaborated. Finally  
32 it is argued that the Wikipedian monitoring of new edits,  
33 especially by its heavy reliance on computational tools,  
34 raises a number of moral questions that need to be  
35 answered urgently.

**Keywords** Bots · Open-source software · Reputation · 37  
Trust · Vandalism · Wikipedia 38

**Introduction** 39

Open-content communities (OCCs) which thrive on con- 40  
tributions from ‘the crowds’ (whether source code, text, 41  
pictures, or videos) have been with us for over two decades 42  
now. Arguably, this movement started with open-source 43  
software (OSS) projects, and further widened with mile- 44  
stone initiatives like Digg, NowPublic, YouTube, and 45  
Wikipedia. The basic parameters of such communities are, 46  
I suggest, twofold. On the one hand they can be distin- 47  
guished by the *goals* they are trying to achieve. Dutton 48  
(2008), however, has eloquently argued that such com- 49  
munities cannot easily be classified by purpose, since the 50  
problems they are trying to solve are bound to change. He 51  
maintains that, instead, they are best characterized by 52  
features of the *technological design* that is implemented: 53  
sharing (1.0), i.e. enabling the bringing together of various 54  
kinds of information; co-contributing (2.0), i.e. enabling 55  
group communication by means of social networking 56  
applications; and co-creation (3.0), i.e. enabling collabor- 57  
ative work by means of tailored software tools for col- 58  
laborative spaces.<sup>1</sup> On the other hand, these OCCs 59  
invariably have to determine conditions of *admission*. Are 60  
all the people showing up to be accepted as contributors? 61  
Or are specific criteria (say, expertise of a kind) to be 62  
applied as conditions of entry? 63

A1 P. B. de Laat (✉)  
A2 Department of Computer Science, University of Groningen,  
A3 Groningen, The Netherlands  
A4 e-mail: p.b.de.laat@rug.nl

<sup>1</sup> Note that West et al. (2012)—an important reference later on in the 1FL01  
section on reputation—fails to make this distinction and lumps all 1FL02  
OCCs together under one label: Collaborative Web Applications 1FL03  
(CWAs). 1FL04

64 Actual communities may position themselves anywhere  
65 along these two parameters. Digg, e.g., is a fully open  
66 ‘social news’ site that enables collective discussion of news  
67 items (co-contributing 2.0). The Encyclopedia of Earth, on  
68 the other hand, is a co-creative (3.0) encyclopedia only  
69 accessible to acknowledged experts. From now onwards in  
70 this article I focus exclusively on those communities that  
71 have adopted the most ambitious approach on *both* counts:  
72 they focus on processes of co-creation with anybody wel-  
73 come. Registration may be obligatory for full participation,  
74 but no criteria for inclusion apply.<sup>2</sup> As explored in a former  
75 study (de Laat 2010), the prime examples of this bold  
76 approach are to be found in source code communities on the  
77 one hand, and textual/pictorial communities on the  
78 other. Well-known examples of the former that take OSS as  
79 their mode of production are Linux, Mozilla, Apache,  
80 and—to mention a more recent one—Drupal. Examples of  
81 the latter with their open wiki spaces are Wikipedia, Citi-  
82 zendum, and Scholarpedia (encyclopedias), and Wikinews  
83 (a citizen journal).<sup>3</sup>

84 The most acute problem these communities have to face  
85 is one of governance: how to manage the incoming flow of  
86 contributions? How to judge the various inputs and process  
87 them? Suppose a range of permissions to perform activities  
88 on project resources has been distinguished: how are these  
89 permissions to be distributed among the crowds? In other  
90 words, which levels of access are to be distinguished (read-  
91 access, write-access, and the like) and to whom are the  
92 distinguished permissions to be distributed?<sup>4</sup> As can  
93 readily be seen, the central issue underlying these choices  
94 is the matter of *trust*: to what extent can potential participants  
95 be trusted to contribute with good intentions and in pos-  
96 session of adequate capabilities? Note that I am not refer-  
97 ring here to personal trust, but to *institutional* trust: the  
98 extent to which the institution approaches its members in a  
99 trusting fashion.

2FL01 <sup>2</sup> Andrea Forte and Cliff Lampe introduce the category of ‘open  
2FL02 collaboration projects’ in their introductory piece to a recent special  
2FL03 issue of the American Behavioral Scientist about such projects (Forte  
2FL04 and Lampe 2013). In my terminology this refers to both co-  
2FL05 contributing (2.0) and co-creation (3.0) communities without barriers  
2FL06 to entry.

3FL01 <sup>3</sup> As a rule, co-created outcomes are licensed to the public with a so-  
3FL02 called Creative Commons licence. This aspect, though of crucial  
3FL03 importance, does not figure in this article and therefore deserves to be  
3FL04 mentioned at least in this footnote.

4FL01 <sup>4</sup> Throughout this article I employ the terms entry/admission and  
4FL02 access as follows. ‘Entry’ or ‘admission’ refers to being accepted as a  
4FL03 participant in the co-creative process; ‘access’ refers to subsequently  
4FL04 obtaining permission to carry out various activities associated with the  
4FL05 process. Compare the—albeit imperfect—analogy of entering a  
4FL06 building through the front door and reaching a hallway (entry,  
4FL07 admission), and subsequently gaining access to the various floors  
4FL08 (access).

## Institutional trust

101 This issue of (institutional) trust can be handled in four  
102 basic ways which can be rendered in concise form as fol-  
103 lows (cf. de Laat 2012c). First, contributors can simply be  
104 *assumed* to be trustworthy, in both moral and epistemo-  
105 logical terms. Without any evidence to warrant the  
106 assumption, contributors are just supposed to be willing to  
107 contribute in an honest and competent fashion. A rationale  
108 is not lacking though; by acting *as if* the other can be  
109 trusted, that other may well turn out to respond in a trusting  
110 fashion. Such ‘normative pressure’ can be conceptualized  
111 in more ways than one. But let me just mention the—as it  
112 would seem—most adequate mechanism to fit our case of  
113 OCCs (based on McGeer 2008). Participants are chal-  
114 lenged to show their capacities as able coders or authors,  
115 and develop them further in the process. Investing in the  
116 capabilities of others may generate its own rewards pre-  
117 cisely by the display of trust involved.

118 A second approach, obviously, is to try and *infer* trust-  
119 worthiness of potential participants. One is on the lookout  
120 for suitable indicators of the kind, such as individual  
121 characteristics, membership of a relevant culture, or links  
122 with respectable institutions. A good reputation may  
123 function likewise. Trust may also be inferred from an  
124 estimation of the costs and benefits inherent in the particular  
125 context. As argued before (de Laat 2010), I do not believe  
126 that OCCs can find many reliable indicators of the kind; a  
127 virtual environment can only yield indicators too fuzzy to  
128 be relied on. The only exception seems to be one’s repu-  
129 tation: as far as a reputation can be established effectively  
130 in cyberspace, it would seem to possess the continuity that  
131 warrants stable inference of a kind.

132 A third mechanism, as recently introduced in de Laat  
133 (2012c), is *backgrounding* trust, which consists of insti-  
134 tuting corrective mechanisms that silently operate in the  
135 background of the community (hence this denomination).  
136 Such backgrounding is comprised of the staging of inten-  
137 sified quality control schemes, especially those that focus  
138 on correcting low-quality contributions of content and/or  
139 actions within the community. These ‘collective monitor-  
140 ing’ schemes are invariably underscored by an etiquette for  
141 proper behaviour which goes beyond the usual legal terms  
142 of use. The norms involved are enforced by sanctioning  
143 deviating members in proper ways, ranging from rebuke to  
144 expulsion. Monitoring schemes and etiquette operating  
145 together are ever so many guarantees that full institutional  
146 trust in community members is warranted.<sup>5</sup>

5FL01 <sup>5</sup> I do not come back anymore to the topic of ‘netiquette’. It is  
5FL02 assumed, by default, that it exists in some form—and is actively  
5FL03 maintained and ‘applied’—in textual/pictorial OCCs. Similarly,  
5FL04 without mentioning it explicitly, it is assumed, as the default again,  
5FL05 that OSS communities are culturally ‘embedded’ in a hacker ethic.

147 The preceding three mechanisms of trust carry the same  
 148 institutional implications for OCCs: unconditional per-  
 149 mission is granted to perform the site activities under  
 150 consideration (such as read- or write-access). Such con-  
 151 clusions no longer apply when a fourth mechanism of  
 152 handling trust is applied: the *substitution* of trust (my  
 153 designation; cf. de Laat 2010). Usually rules and regula-  
 154 tions are introduced in evolving communities in order to  
 155 manage the interactions between members. Often, these  
 156 exhibit characteristics of encroaching bureaucracy:  
 157 hierarchical distinctions and vertical control appear on the  
 158 scene. As a result, participants' discretion to perform tasks  
 159 becomes circumscribed. Thereby, unwittingly or not, the  
 160 problem of institutional trust is tackled by the (partial)  
 161 elimination of the need for granting trust; trust is being  
 162 substituted. In the case of our OCCs this governance  
 163 mechanism entails a delineation of the conditions under  
 164 which members may get access to the various activities  
 165 within the community. At a minimum one layer of super-  
 166 vision is introduced. Normally, the role hierarchy is more  
 167 extended; as a result, the blanket granting of immediate  
 168 contributory access is eliminated from the repertoire. Roles  
 169 and permissions have to be earned; they are no longer  
 170 granted indiscriminately.<sup>6</sup>

171 The central research question addressed in this paper is  
 172 the following. How do the communities that focus on fully  
 173 open co-creation (OSS communities on the one hand and  
 174 textual/pictorial communities on the other) manage the  
 175 contents that are contributed? In particular, to what extent  
 176 can they be perceived to rely on each of the four mechanisms  
 177 outlined above? In broad terms it is shown below that the  
 178 OSS communities continue to rely considerably on hier-  
 179 archy (the 'onion' model) as a mode of substituting trust.  
 180 The communities for text/pictures, though, eschew the  
 181 tools of hierarchy to a large extent. Furthermore, as a  
 182 means to fight 'vandalism', Wikipedia in particular has  
 183 developed its own distinct mechanism of backgrounding  
 184 trust: the community is mobilized to monitor incoming  
 185 contents closely. This research charts the associated pro-  
 186 cess of developing new monitoring tools. Moreover,  
 187 approaches in computational science are at the basis of  
 188 even more sophisticated monitoring tools, and have also  
 189 led to the creation of software bots that autonomously scan  
 190 for vandalism. Finally, measuring reputation on a continuous  
 191 basis is being contemplated in Wikipedian circles, since  
 192 that indicator would allow monitoring more closely the  
 193 edits from low-reputation contributors—as presumably  
 194 being the least trustworthy of all. Let it be remarked in

advance, that these findings are summarized in a table at  
 the end for easy reference (Table 1).

**Source code**

Several publications have investigated what this challenge  
 of trust means for the communities that stand at the origin  
 of the open-content movement: OSS (Crowston et al. 2004;  
 Holck and Jørgensen 2005; de Laat 2007, 2010). These  
 communities were found to rely mainly on the mechanism  
 of substitution of trust: *hierarchy* is the standard solution.  
 In what is aptly denominated the 'onion' model (Crowston  
 et al. 2004), several layers are distinguished that obtain an  
 ever increasing number of permissions to perform activities  
 on the site's resources. The number of layers may vary  
 from three upwards. A typical onion (as employed on  
 Tigris) consists of the following three roles. An observer  
 has read-only access to most of the project's documenta-  
 tion and files. In this capacity he (obviously) may return  
 any comments in text or code he wishes to contribute. A  
 developer obtains the additional permission of write-  
 access: inserting code in files of the official tree and/or text  
 in other project files. The project owner at the top (the  
 onion's core) manages the project as a whole (and in this  
 capacity also decides on whether someone is to rise in the  
 ranks or not).

As an OSS project grows in size, there seems to be a  
 noticeable tendency to expand this hierarchy and define  
 ever more roles and associated conditions. Take Mozilla,  
 with their 80 modules definitely a larger project (more  
 details in de Laat 2010). For one thing, additional roles like  
 'super-reviewer', 'release driver', and 'benevolent dictator'  
 (for resolving conflicts) have evolved. For another, write-  
 access concerning code is no longer a straightforward  
 permission. After testing proposed code in their own copy  
 of the official tree developers are urged, before actually  
 committing to the official tree, to ask for a twofold per-  
 mission: from the owner of the specific module they hap-  
 pen to work in, and from one of the 'super-reviewers' who  
 guard the quality and consistency of the overall Mozilla  
 code base.

These hierarchical solutions seem to have been  
 employed for at least two decades now. They can be  
 considered stable and robust solutions to the problem of  
 trust (more narrowly) as well as coordination (more gen-  
 erally). This conclusion is unaltered by the recent devel-  
 opment of newer kinds of versioning systems. The original  
 ones (like CVS) were centralized in client-server fashion,  
 with all traffic directed to and from one canonical code  
 base. The newer ones (like Bitkeeper, Git, and Mercurial)  
 are distributed systems. Each participant can obtain an  
 integral copy of the public main repository, with all history

6FL01 <sup>6</sup> Notice that the essential differences between the third and fourth  
 6FL02 mechanisms of trust management are fleshed out and explored more  
 6FL03 extensively in the sections that follow, culminating in a more  
 6FL04 elaborate analysis under the section 'Collective monitoring within  
 6FL05 Wikipedia: interpretation'.

245 included—his own public fork. He can then experiment on  
 246 his own in a private copy of this, and finally publish his  
 247 code changes back into his public repository. So each and  
 248 every one publicly shows off his own fork, with his own  
 249 code improvements incorporated. It is then up to each  
 250 participant to ‘pull’ any commits from others into his own  
 251 private copy, and if found to be satisfactory, to ‘push’ it  
 252 into his public fork (cf. the clear exposition on [http://eqqon.com/index.php/Collaborative\\_Github\\_Workflow](http://eqqon.com/index.php/Collaborative_Github_Workflow)).<sup>7</sup>

253 Some herald the new system as the pinnacle of peer-to-  
 254 peer production, implying that finally all participants may  
 255 become full committers and are therefore cooperating on an  
 256 equal basis (Orsila et al. 2009). This would seem to be an  
 257 overblown interpretation. What happens is, that participants  
 258 are more able to ‘advertise’ their code changes (and if need  
 259 be, to steer their public fork in another direction than the  
 260 main project, thereby effectuating what is normally under-  
 261 stood by ‘forking’). However, it is still up to the owners and  
 262 developers of the *original* project—or any other for that  
 263 matter—whether they actually incorporate any changes into  
 264 *their* public repository (by pulling and pushing commits). In  
 265 all forks involved, particularly in the main one (the ‘official’  
 266 project), hierarchy continues to obtain.

267 Rising in this hierarchy is simply a matter of proving  
 268 oneself able enough (more details in de Laat 2010). For the  
 269 purpose of inferring trustworthiness as concerns fulfilling  
 270 higher roles three indicators are being used. Most important  
 271 are hacking skills. Only after providing some excellent con-  
 272 tributions may one successfully apply for developer status. In  
 273 order to rise higher, proof of leadership skills is also required.  
 274 Of late, due to rising concerns about sloppy, buggy, or Trojan-  
 275 horse code, in some projects ideological allegiance to the  
 276 cause of OSS has also become a requirement (e.g., Debian).

277 With OSS, therefore, the issue of trust has found a stable  
 278 solution: a division of roles is employed (*substitution* of  
 279 trust), the ranks of which are filled in accordance with  
 280 proven achievements inside the project (*inference* of trust)  
 281 (cf. Table 1). Some have phrased this combination: a ‘role  
 282 meritocracy’. In such a model, trust is not granted at the  
 283 outset. Only as one’s participation develops in satisfactory  
 284 fashion, ever more permissions may be forthcoming. In  
 285 that sense, trust is not granted *ex ante*, but *ex post*, step by  
 286 step, to the extent that one has proven oneself to be an able  
 287 and loyal hacker.

## 289 Text and pictures

290 Fully open co-creative communities other than OSS started  
 291 only a decade ago. Inspired by the successful approach of

292 producing source code in collaborative fashion, the pro-  
 293 duction model migrated from source code to content con-  
 294 sisting of text and pictures. If we omit the attempts at  
 295 writing wikibooks as being insignificant,<sup>8</sup> the field is  
 296 dominated by projects that focus on producing encyclo-  
 297 pedias (Wikipedia, started in 2001; Scholarpedia, 2006;  
 298 Citizendium, 2007), and journals (Wikinews, 2004).<sup>9</sup> As  
 299 stated, everybody without exception is welcome to  
 300 participate in co-creation. Almost all of these projects,  
 301 however, take the trusting approach to the next level—as  
 302 never contemplated in OSS circles: not only may every-  
 303 body contribute; one’s contributions are also ‘published’  
 304 right away, visible to anyone. Since all projects employ the  
 305 wiki tool, this means that, next to the obvious read-access,  
 306 full and immediate *write-access* to the wikis is granted.  
 307 ‘Real time’ contributing is the motto.<sup>10</sup> Notice moreover,  
 308 that the look-and-feel of the websites’ architecture is very  
 309 similar (while directly modelled after Wikipedia which was  
 310 one of the first sites to use wikis).

311 Right at the outset it should be emphasized, however, that  
 312 the way in which the wiki tool is applied, sets Scholarpedia  
 313 apart from all others. This natural sciences project operates  
 314 more like a scientific journal than as the Wikipedia we all  
 315 know (details to follow obtained from [www.scholarpedia.org](http://www.scholarpedia.org)). The major part of production is carried out in a *closed*  
 316 wiki space. Experts are invited to become the author of a  
 317 particular entry. After it has been written, it has to be  
 318 reviewed by invited experts: only after their approval does  
 319 the entry (signed by its authors) become publicly accessible  
 320 to all. This whole process leading up to publication is a  
 321 hierarchical one. Several layers can be distinguished with  
 322 increasing powers: ‘scholars’ may review, and ‘editors’ at  
 323 the top take care of overall coordination.

324 When the wiki space of the article opens up to the public  
 325 for comments, ‘curators’ are appointed that carry  
 326 responsibility for the article’s further evolution. This  
 327

8 Attempts to write a book collectively (‘networked book’) failed. In 8FL01  
 this vein several projects were initiated, of which the most famous one 8FL02  
 was staged by Penguin (dubbed ‘A Million Penguins’), inviting the 8FL03  
 crowds to produce a book together (2007). In a time span of 5 weeks a 8FL04  
 ‘wikinovel’ was produced, with some 1,500 people contributing 8FL05  
 (Pullinger 2012). From our perspective (of trust) the main observation 8FL06  
 to be made is that, due to the many reactions ultimately verging on 8FL07  
 vandalism, a team of students had to practise filtering of incoming 8FL08  
 edits – a hierarchical kind of arrangement that figures later in the main 8FL09  
 text as well. 8FL10

9 Note that more such general encyclopedias have actually been 9FL01  
 initiated during the last decade, many of them copying the software of 9FL02  
 the Wikipedian production model (available as ‘open source’). They 9FL03  
 are not taken into consideration here since they are either only part- 9FL04  
 encyclopedias, or carry a distinct ideological message, or have simply 9FL05  
 not survived. 9FL06

10 The exceptions to this rule either slightly qualify full write-access 10FL01  
 (Wikinews’ front page), or never introduced it in the first place 10FL02  
 (Scholarpedia); to be commented on below. 10FL03

7FL01 <sup>7</sup> I felt free to use the masculine personal pronoun in the paragraphs  
 7FL02 above, since almost all developers concerned are male.



328 curator, typically one of the authors who wrote the article  
 329 in the first phase, is bestowed with hierarchical powers as  
 330 well: edits that come in need to obtain his/her approval  
 331 before appearing in the official wiki version. In common  
 332 parlance: they *filter* edits for acceptance. A requirement for  
 333 contributing to the wiki, finally, is registration as a user. In  
 334 accordance with the ‘assurance view’ as set forth by Moran  
 335 (2005), the act of registration may be considered a sign of  
 336 trustworthiness. The contributor declares to stand behind  
 337 his/her (future) edits and assumes responsibility for them.  
 338 Some guarantee for their truth is provided. As a result, the  
 339 relationship between contributor and community is turned  
 340 into a normative one.

341 Summarizing, the trusting invitation of write-access that  
 342 Scholarpedia extends to the general public is relatively  
 343 small: it only applies to the last stage of ‘refinement’ of  
 344 articles. The encyclopedia manages the inherent problem  
 345 of trust by letting curators filter edits (*substitution* of trust),  
 346 and by asking users to register (*inference* of trust) (cf.  
 347 Table 1).

348 At Citizendium, Wikinews, and Wikipedia the gesture  
 349 of trust towards the general public is much broader: the  
 350 logic of full read- and write-access to the wiki is applied  
 351 *throughout* the production process. How is the problem of  
 352 trust in fully fledged form handled by these communities?  
 353 At Citizendium, by far the smallest of them, the means of  
 354 governance are quite slim. On the one hand (as in Schol-  
 355 arpedia), every contributor has to register (by ‘real name’).  
 356 On the other hand, constables (similar to administrators in  
 357 Wikipedia, cf. below) are appointed to act as policemen  
 358 when interactions derail and resolution by brute force  
 359 seems the only option. These minimal means to handle the  
 360 trust problem—instances of, respectively, the *inference* and  
 361 the *substitution* of trust (Table 1)—suffice until now to  
 362 streamline interactions within the Citizendium community.

363 As far as Wikinews is concerned, their governance is quite  
 364 slim as well. They are actually very similar to Wikipedia—  
 365 not surprising, since both fall under the umbrella of the  
 366 Wikimedia Foundation. Next to the appointment of admin-  
 367 istrators, in 2008 Wikinews (at least the English version)  
 368 introduced a reviewing system for the last phase of produc-  
 369 tion: any article from the wiki ‘newsroom’ has to be  
 370 approved of before appearing on the ‘main page’. These tools  
 371 of governance—both instances of the *substitution* of trust  
 372 (Table 1)—are elucidated below.

373 Wikipedia: early governance

374 The means of governance are no longer minimal, however, for  
 375 Wikipedia, with traffic a thousand times larger (in terms of  
 376 articles, users, their edits, and the like). In particular, it gets  
 377 confronted with vandalism on a large scale: current estimates  
 378 hover around nine thousand vandalist edits (7 % of all edits)

daily.<sup>11</sup> Full write-access, on that scale, is no longer an easy 379  
 undertaking.<sup>12</sup> 380

381 In the old times, when they started, Wikipedians may  
 382 have thought they could get by just by appointing  
 383 ‘administrators’ who have the powers to protect and delete  
 384 pages, and block users (either temporarily or permanently).  
 385 These are appointed by higher ‘bureaucrats’. This hierarchical  
 386 arrangement—an instance of *substituting* trust (Table 1)—  
 387 was supposed to take care of disruptive behaviours on the  
 388 site.

389 Soon enough, however, it became clear that possible  
 390 disruptions in the fully open access condition (without  
 391 registration requirement, implying that anonymous editing  
 392 is possible) could not be held in check with this minimal  
 393 hierarchy. In response, many initiatives have been  
 394 unfolding over the years. Early onwards, additional per-  
 395 missions (or flags) were developed that constituted ever  
 396 more technical tools to get to grips with disruptive con-  
 397 tributors. The ‘rollback’ permission allows to use a spe-  
 398 cially installed button that quickly reverts consecutive bad  
 399 edits on a page by one and the same user; the ‘checkuser’  
 400 permission allows to see all IP-addresses as used by a  
 401 supposed ‘vandal’; and the ‘oversight’ permission enables  
 402 suppressing edits and make them disappear (almost)  
 403 completely from the files (used for materials that are  
 404 defamatory, intrude privacy, or violate copyright) (for all  
 405 permissions cf. WP:UAL). Moreover, as a means of  
 406 intrusion prevention, the ‘abuse filter’ permission allows  
 407 to write and install automatic filters on incoming traffic.  
 408 As to be expected, these permissions were only granted  
 409 on a minimal basis: to (a selection of) administrators,  
 410 buttressing the hierarchy. The rollback permission in  
 411 particular was also granted to some more registered users  
 412 who had shown in practice that they could be trusted.

413 When the problem of vandalism persisted, Wikipedia  
 414 entered a phase in which a review system was contem-  
 415 plated (full details in de Laat 2012b). Incoming edits were  
 416 to be reviewed for evidence of vandalism *before* being  
 417 accepted and appearing on the screen (‘official version’).

11 These observations refer to the *English* version of Wikipedia. In 11FL01  
 the remainder of this article, unless specified otherwise, I always refer 11FL02  
 to that language version—actually the largest of all language versions 11FL03  
 of Wikipedia. 11FL04

12 Some conditions have been introduced in Wikipedia that qualify 12FL01  
 write-access for all (WP:UAL). Any user, without an account 12FL02  
 (‘unregistered’), may read and edit entries (pages). Upon registration, 12FL03  
 the user may *also* create new pages. After some time (four days and at 12FL04  
 least ten edits) the registered user automatically becomes ‘autocon- 12FL05  
 firmed’, which implies that (s)he may *also* move pages around and 12FL06  
 upload files and pictures. Write-access may be said to be ‘complete’ 12FL07  
 by then. Currently over a million (English) users are autoconfirmed. 12FL08  
 Let me remark finally, in order to avoid any misunderstanding, that 12FL09  
 write-access does not only involve the right to add or change text but 12FL10  
 also to *delete* text. 12FL11

418 Depending on the parameters chosen, the system can  
 419 assume various shapes. Let me just elucidate the system as  
 420 adopted for a yearlong trial (2010–2011) in the English  
 421 Wikipedia (known as ‘flagged revisions’). Edits to specific  
 422 (sensitive) entries were put on hold (‘pending changes’).  
 423 As soon as versions with new edits were approved, they  
 424 were flagged (‘flagged revision’) and promoted to be the  
 425 official version. The reviewers involved had to apply to the  
 426 administrators and show over a hundred accepted edits and  
 427 an impeccable track record as far as vandalism and  
 428 harassment is concerned. As can be seen, this represented a  
 429 further incursion into bureaucratic terrain, with a new layer  
 430 of reviewers in operation. This time the full write-access  
 431 permission to users came under siege (at least in those  
 432 spaces where the experiment applied): ordinary users  
 433 across the board came under scrutiny from more experi-  
 434 enced users who engaged in filtering their edits.<sup>13</sup> In my  
 435 terminology: one more step towards the *substitution* of  
 436 trust was under consideration (Table 1). No wonder that  
 437 the reviewing system was heavily condemned by many as  
 438 ‘just added bureaucracy’ and had to be abandoned after the  
 439 trial period. As of now, only entries subjected to acute  
 440 quarrels (often those about living persons) can be brought  
 441 under this flagging regime.<sup>14</sup>

#### 442 Wikipedia: collective monitoring

443 Instead of filtering, Wikipedia now mainly relies on  
 444 another approach to confront vandalism. It involves a  
 445 permanent mobilization of Wikipedians to fight low-quality  
 446 contributions and their authors, and keep vandalism at  
 447 bay. Over and beyond any normal interactions in the wiki  
 448 spaces, initiatives are unfolding to weed out vandalism and  
 449 disruptions. This campaign of close watch proceeds more  
 450 silently in the background; hence my denomination of the  
 451 mechanism involved as *backgrounding* trust (Table 1).  
 452 This vigilance is several years old now; increasingly,  
 453 software tools are being developed to support this cam-  
 454 paign. Let me explain.<sup>15</sup>

455 First, individuals and groups are called upon to organize  
 456 themselves and be alert to vandalism. The main focus is on

new edits, whether large or small.<sup>16</sup> These can be displayed  
 in ‘real time’ by using the Lupin tool which also allows  
 making a selection of them (such as edits containing  
 ‘suspect words’, ignoring administrator edits or talk pages)  
 (WP:Lupin). In a more tailored fashion users can maintain  
 their own personal ‘watch lists’; each entry on such a list is  
 kept under constant watch for new edits that come in.  
 Subsequently the ‘patroller’ has to make the decision  
 whether to accept an edit under scrutiny or delete it as an  
 instance of vandalism. In order to facilitate that process  
 various useful buttons can be installed (on the patroller’s  
 page). A tool such as Twinkle installs buttons for various  
 kinds of rollback (qualifying the edit as vandalism, dubi-  
 ous, or made in good faith) and for easily accessing the talk  
 page of a detected vandal user and attaching a warning  
 template (WP:Twinkle).

The steps of detection and action are nicely combined in  
 the integrated tool called Huggle (WP:Huggle). It allows  
 displaying fresh edits using various filters: all edits; only  
 edits by anonymous users; only edits by users with warn-  
 ings on their talk pages; only edits by humans (bots  
 excluded, see below); and so on. Moreover one may focus  
 on edits that have a high probability of being vandalistic, as  
 determined by an algorithm.<sup>17</sup> Subsequently the patroller  
 may delete the edits identified as vandalist and warn their  
 authors on their respective talk pages. As this is a poten-  
 tially dangerous tool, one needs the rollback permission for  
 it.<sup>18</sup>

In recent years, this monitoring approach has obtained a  
 fresh impulse from several developments in computer  
 science. These all revolve around identifying the quality of  
 (Wikipedian) edits or articles. Let me single out some

13FL01 <sup>13</sup> Note the analogy with the division of roles in OSS: observers’  
 13FL02 contributions have to be scrutinized by developers before acceptance.

14FL01 <sup>14</sup> Based on the likely introduction of this flagged-revisions scheme,  
 14FL02 some time ago I foresaw a convergence of the designs for open-source  
 14FL03 software and encyclopedias (de Laat 2010). It has now become clear  
 14FL04 that this convergence is not taking place.

15FL01 <sup>15</sup> As demonstrated in de Laat (2012c), social news sites and citizen  
 15FL02 journals similarly rely on backgrounding trust: voting schemes push  
 15FL03 high quality articles to a prominent or visible position—and likewise,  
 15FL04 relegate low quality contributions to an inconspicuous or invisible  
 15FL05 position. These sites are not considered here, however, as they are of  
 15FL06 the 2.0 co-contributing category.

16FL01 <sup>16</sup> On a more personal note, let me quote from my own recent  
 16FL02 experience of vandalism patrolling: words may be inserted (yo-  
 16FL03 dolohee, poo, popcorn, peanut butter), substituted (Boeing 747  
 16FL04 Dreamliner is changed into Nightmareliner; a hip hop album sales  
 16FL05 figure is changed from 295,000 to 2,295,000), or whole paragraphs  
 16FL06 blanked (or replaced by HAHAHA). Vandalist insertions can also be  
 16FL07 larger, and even be creative. Let me give the example of the entry  
 16FL08 ‘Heat Pipe’, in the middle of which the following lines were inserted  
 16FL09 (at 20:42 on 28 February 2013): “A little known fact is that number of  
 16FL10 Dwarfs actually live inside these pipes and help with constant  
 16FL11 maintenance, they may need to be replaced at some point in the  
 16FL12 computers life due to wars between the dwarfs that end with  
 16FL13 numerous casualties. Treaties have been implemented between the  
 16FL14 dwarf clans, but they can never live in harmony.”

17FL01 <sup>17</sup> An early example of ‘algorithmic power’, to be discussed more  
 17FL02 fully below.

18FL01 <sup>18</sup> In a larger vein not only new edits but also new *entries (pages)* as a  
 18FL02 whole are watched constantly. ‘New Pages Patrol’ is a system that  
 18FL03 signals newly created pages and invites Wikipedians to check whether  
 18FL04 or not these conform to various criteria (concerning not only  
 18FL05 vandalism, but also relevance, substance, harassment, advertising,  
 18FL06 copyright violations, etc.) (WP:NPP). Unwelcome candidates are to  
 18FL07 be nominated for so-called ‘speedy deletion’. This patrol is intended  
 18FL08 to eradicate quality problems right from the start.

490 specific approaches that relate to the focus on vandalism.  
 491 The ‘Wikitrust model’ (cf. Adler and de Alfaro 2007)  
 492 employs a specific type of metrics to gauge quality (or  
 493 credibility) of entries: the ‘survival’ of individual edits  
 494 during their evolution. The term survival is connected with  
 495 the central metaphor of voting: each round of editing is  
 496 seen as casting a vote on edits in view at the moment. The  
 497 longer the period over which edits remain intact, the more  
 498 they become credible. Associated with this, every time an  
 499 edit survives a vote its author earns an increase in reputa-  
 500 tion. As a result, authors obtain an increasing reputation (as  
 501 trustworthy Wikipedians) the more they edit *and* the edits  
 502 involved ‘survive’. Obviously, the reverse applies as well:  
 503 any edit deletion punishes its author by a decrease in  
 504 reputation. As a more subtle point, not every Wikipedian  
 505 counts equally in the process. For one thing, those of a  
 506 higher reputation, upon endorsing edits, increase their  
 507 credibility to a larger extent than those of a lower reputa-  
 508 tion.<sup>19</sup> For another, a vote cast by Wikipedians of high  
 509 reputation counts more for someone’s reputation than a  
 510 vote by a low reputation Wikipedian. That is to say,  
 511 Wikipedians of high reputa have both more credibility  
 512 points and more reputational points at their disposal to  
 513 distribute than Wikipedians of low reputation.<sup>20</sup>

514 Furthermore, computational approaches to detect van-  
 515 dalist edits have been worked on extensively. As of now,  
 516 they can be classified in four categories (Adler et al. 2011).  
 517 Each has its own focus. First, features of language can be  
 518 inspected (e.g., bad words, pronoun frequency). Secondly,  
 519 textual features (language-independent) can be the focus  
 520 (e.g., use of capitals, changes to numerical content, dele-  
 521 tion of text). Thirdly, metadata of edits can be indicators of  
 522 vandalism (e.g., anonymity, local time the edit was made,  
 523 absence of revision comment). Finally, the measure of an  
 524 editor’s reputation—as elaborated above—may be useful: a  
 525 low reputation makes vandalism more likely. A measure of  
 526 country reputation is also in use. It is obviously a challenge  
 527 for computer scientists to determine which type of  
 528 approach yields the best vandalism detection scores.  
 529 Recent experiments indicate that a *combination* of all four  
 530 may deliver the best results (Adler et al. 2011).

531 These approaches are used to develop practical tools:  
 532 autonomous software bots for vandalism detection and  
 533 repair. With overall hundreds of bots having been

534 developed by the Wikipedian community, several of them  
 535 have a specific focus on vandalism and are currently  
 536 operative. Until a few years ago, most of them were based  
 537 on detection of suspect linguistic or textual features (e.g.,  
 538 ClueBot). They intervened automatically when suspicious  
 539 words (enumerated on ‘black lists’) were inserted or whole  
 540 pages were blanked: the edit was reverted and a note of  
 541 warning placed on the suspect’s talk page. Remarkably, a  
 542 newer generation of bots takes a quite different approach to  
 543 vandalism detection by operating like a neural network.  
 544 The bot gradually learns to distinguish bona fide edits from  
 545 vandalist edits. For the purpose, it has to be ‘fed’ with real  
 546 examples of both kinds of edits. A critical feature is its  
 547 false positive rate: it is set at or just below the rate that  
 548 ordinary humans achieve. The successor to ClueBot,  
 549 ClueBotNG, operates like that. All such bots are allowed to  
 550 scan Wikipedian spaces (using a ‘bot account’), but only  
 551 after heavy testing, public discussion, and permission from  
 552 the ‘Bot Approvals Group’. Notice finally that they heavily  
 553 contribute to vandalism reversal: the top scorers among  
 554 them, whether from the old or the new generation, have  
 555 performed millions of edits each. That is more than  
 556 ordinary humans can ever hope to achieve.

557 In a final step, any of the four vandalism detection  
 558 algorithms can be built into integrated tools for both edit  
 559 detection and action. As a result, such tools become more  
 560 powerful. The promise is that the combination of auto-  
 561 mated and human power will yield better results than each  
 562 on their own. A prime example of such ‘assisted editing’ is  
 563 the STiki tool (WP:STiki). At the back-end (processing of  
 564 edits), fresh edits are continuously monitored for vandalism,  
 565 based on the metadata approach (cf. above; a ClueBotNG  
 566 engine is also built in now). At the front-end (the GUI),  
 567 operators get to see the top edit of a queue of suspect edits,  
 568 ordered by (presumed) vandalism scores. In response, edits  
 569 can either be accepted (classified as ‘pass’ or ‘innocent’),  
 570 or be reverted as unconstructive (classified as either ‘in bad  
 571 faith’ or ‘in good faith’) and their authors be given a  
 572 warning. This tool is far superior to the Huggle tool. For  
 573 one thing, detection has fully become algorithm based; for  
 574 another, edits under review by someone are ‘reserved’ (no  
 575 simultaneous checking), and ‘innocent’ edits leave the  
 576 queue and are therefore not re-inspected. As it is so  
 577 powerful and may easily wreak havoc on wikispaces, the  
 578 tool requires special permission and is usually only granted  
 579 to users with rollback permission (or similar status).

580 For the future mixed models are also being contem-  
 581 plated, as a kind of mid-solution in between fully auto-  
 582 mated bots and humans armed with assisted editing tools.  
 583 In them, while vandalism detection is automated based on  
 584 one or more algorithms, subsequent action is both human-  
 585 based and computer-based. For example, imagine the fol-  
 586 lowing system (which combines separate elements as

19FL01 <sup>19</sup> In actual fact, as soon as someone’s reputation is too low in  
 19FL02 relationship to the credibility of a specific edit, endorsing the edit does  
 19FL03 not increase its credibility at all.

20FL01 <sup>20</sup> This model has resulted in a practical tool: the WikiTrust  
 20FL02 extension (Adler et al. 2008). It continuously calculates the credibility  
 20FL03 of words in an entry as ‘voting’ continues and assigns colours to them  
 20FL04 accordingly (ever lighter shades of orange indicate old age). The tool  
 20FL05 may assist users to focus their efforts on the fresh parts of the text  
 20FL06 (dark orange).



587 mentioned in West 2011). Incoming edits are first sorted by  
 588 their vandalism scores. Subsequent action for reversal then  
 589 depends on that score. Edits with high probability of van-  
 590 dalism are rejected automatically (without ever appearing  
 591 on the screen); all other edits do get accepted and become  
 592 part of Wikipedia. Subsequently, though, edits with  
 593 medium probability of vandalism among them (considered  
 594 to be suspect) are suggested to human patrollers for making  
 595 a decision. As a result, the decision-making process char-  
 596 acteristic of STiki becomes, as it were, ever more auto-  
 597 mated and autonomous.<sup>21</sup>

## 598 Wikipedia: reputation tracking

599 A final subtlety needs to be described concerning the col-  
 600 lective monitoring efforts in Wikipedia. In the future these  
 601 may become intertwined with a specific type of indicator  
 602 for inferring trustworthiness: reputation. That indicator  
 603 might then influence the amount of monitoring deemed to  
 604 be necessary: high reputation would render monitoring  
 605 superfluous, while low reputation would necessitate an  
 606 increase in monitoring. Monitoring becomes *differentiated*  
 607 along the dimension of (imputed) reputation. Let me  
 608 explain.

609 Many OCCs keep track of a contributor's reputation  
 610 within their particular community. Often denoted as  
 611 'karma', it provides a judgment about one's achievements  
 612 condensed into a single numerical score. Usually, the  
 613 measure selected for the purpose is quite simple and intu-  
 614 itive. In social news sites up votes (+1) and down votes  
 615 (-1) on one's contributions (modelled after 'digging' and  
 616 'burying' as pioneered by Digg) are added up to produce  
 617 one's karma. In many citizen journals, one obtains points  
 618 from one's various types of contributions, and from the  
 619 comments they evoke from others in return; the sum total  
 620 of these points is considered an indicator of reputation. In  
 621 Wikipedia itself, the total number of ('reviewed') edits one  
 622 has contributed is an accepted measure of reputation.

623 Now, for what purposes do OCCs keep track of their  
 624 members' reputation—for whatever it is worth? The main  
 625 intent seems to be motivating members to continue the  
 626 good work within the community. For this purpose, repu-  
 627 tational scores are displayed publicly on user pages (either  
 628 more discreetly, or more prominently). In addition, some  
 629 communities compose 'leader boards' and 'member  
 630 rankings' from the reputational scores and display them on

a highly visible spot. Similarly, 'recognition awards' 631  
 (Ground Report) and 'barn stars' (Wikipedia) are awarded 632  
 to prolific members. The phenomenon that we observe here 633  
 is the 'gamification' of community work: the introduction 634  
 of game design elements in the non-game environment of 635  
 these communities (as the classic definition of gamification 636  
 is usually formulated: Deterding et al. 2011).<sup>22</sup> 637

In a few communities more innovative use is made of 638  
 reputation, in an effort to realize the potential of that 639  
 measure for justifying the distribution of 'privileges'. 640  
 These do not relate to editing as such, but most often to 641  
*control over* editing. In Slashdot, only members of good 642  
 reputation can be invited by site editors to assist with the 643  
 task of moderation (rating articles as either constructive or 644  
 not).<sup>23</sup> In HackerNews, similarly, high karma members 645  
 obtain the privilege to flag items as abusive (for subsequent 646  
 verification and action by the editorial team). In Wikipedia 647  
 itself, finally, the same kind of reasoning has led to the 648  
 requirement of a minimum edit count for anyone volun- 649  
 teering to become an official 'reviewer' (in the English 650  
 'flagged-revisions' scheme the norm is one hundred edits, 651  
 in the German 'gesichtete-Versionen' scheme it is three 652  
 hundred edits). 653

But reputation can also be conceived as useful for the 654  
 background process *itself* of monitoring new edits coming 655  
 in as just described. The hunch is that the lower a con- 656  
 tributor's reputation, the less (s)he can be trusted to be a 657  
 good Wikipedian; accordingly, his/her edits are to be 658  
 watched closely. For this purpose, simple edit count (as 659  
 mentioned above) is too raw a measure, as it can hardly be 660  
 interpreted as an indicator of quality contributions. The 661  
 Wikitrust model (cf. above) meets these concerns: it pro- 662  
 poses a far more sophisticated measure of reputation (sum 663  
 of edits that effectively survived the process of collective 664  
 'voting'). This measure changes dynamically up and down 665  
 in accordance with how contributors' edits evolve. Pre- 666  
 cisely for that reason, this kind of reputation is the cor- 667  
 nerstone of one of the main algorithms of vandalism 668  
 detection (the fourth one, as elucidated above). As such it 669  
 can be incorporated as a detection engine in any of the 670  
 integrated 'assisted editing' tools. In STiki, e.g., it had been 671  
 integrated as one of four engines in the back-end, enabling 672  
 human operators to choose the reputational queue of edits 673  
 for inspection and focus their vandalism detection efforts 674  
 accordingly.<sup>24</sup> 675

21FL01 <sup>21</sup> A similar experience can be obtained from using the experimental  
 21FL02 tool wpcvn.com. It presents possible instances of vandalism that  
 21FL03 occurred over the last hour to its human operators, combined with the  
 21FL04 'karma' (i.e., reputation) of their authors. Only edits performed by  
 21FL05 contributors with *negative* karma are shown. The design thus steers  
 21FL06 attention to the low-karma-contributors. This tool, however, is no  
 21FL07 longer working as of January 2014.

<sup>22</sup> Note that the number of edits to Wikipedia patrolled by means of 22FL01  
 STiki is also kept track of on a 'leader board' – 'assisted editing' itself 22FL02  
 is also subjected to gamification. 22FL03

<sup>23</sup> The privilege rotates regularly over the Slashdot population-of- 23FL01  
 high-repute as a whole, thereby avoiding role fixation. 23FL02

<sup>24</sup> Due to operational difficulties this reputational type of engine for 24FL01  
 STiki is now out of order. 24FL02



676 In such and similar instances, the model of *backgrounding*  
 677 trust within Wikipedia becomes intertwined with *inferring*  
 678 trustworthiness from the indicator of reputation (e.g., as  
 679 following from the Wikitrust model) (included in Table 1).  
 680 This indicator is supposed to *optimize* monitoring efforts.  
 681 Nevertheless, such use of reputation, whether in the narrow  
 682 sense (vandalism detection) or in the broader sense (distrib-  
 683 ution of roles and privileges), is a controversial issue at the  
 684 moment. This has to do with three main problems.

685 For one thing, it proves very difficult to construct a  
 686 satisfying operational measure of reputation (the following  
 687 is largely based on West et al. 2012, and Adler and de  
 688 Alfaro 2007). Ideally, the measure should rise and fall,  
 689 reflecting increases and decreases in imputed trustworthi-  
 690 ness accurately. To that end, with a range say from 0 to 1,  
 691 its starting value should lie somewhere in the middle (1/2).  
 692 That would reflect a neutral evaluation of newcomers. With  
 693 such a midway starting value, however, vandals (whose  
 694 reputation will plummet to 0 after a series of vandalist  
 695 actions) always dispose of the option to open a new  
 696 account and thereby start all over again with an unblem-  
 697 ished reputation (a mechanism dubbed ‘karma bankruptcy’  
 698 by Farmer and Glass 2010: 161–162). In other words, such  
 699 a measure would not provide incentives to abstain from  
 700 vandalism—it bears no cost. Compare a reputation in eBay:  
 701 when it has become too low, sellers simply open a new  
 702 account in order to continue selling. To avoid this  
 703 unwanted mechanism, the reputational start value should  
 704 be at the bottom of its range (close to zero)—notice that the  
 705 Wikitrust measure has this characteristic. Then, however,  
 706 another problem surfaces: both starters and vandals have  
 707 the same reputation and receive the same treatment (such  
 708 as being put under close watch)—also a clearly undesirable  
 709 feature from the community point of view.

710 For another, the question whether reputation acquired  
 711 should be made public or not is a difficult one (again, cf.  
 712 West et al. 2012). In order to function as an incentive for  
 713 proper behaviour, public visibility is a requirement. Or, as  
 714 a variety, a community should have a clearly stated policy  
 715 that reputations are kept for purposes of governance—  
 716 without necessarily making them known. After all, partic-  
 717 ipants only need to know *that* a reputational mechanism  
 718 exists which influences their fortunes in the community.  
 719 Such public awareness, however, creates incentives to  
 720 optimize one’s reputation—in ways that are not necessarily  
 721 to the benefit of the community. Inside Wikipedia one such  
 722 ‘gaming’ tactic would be to contribute a series of small  
 723 edits in succession, instead of the whole text as one single  
 724 large edit. But then, consider the alternative option of  
 725 keeping track of reputations in secrecy. This would allow  
 726 governance-by-reputation to continue all the same.

727 However, it would take away the incentive to behave  
 728 well—and could be considered rather sneaky at that.

729 In addition, the computations involved are extremely  
 730 complex and consume large resources. The reputation of  
 731 the contributors involved would have to be updated with  
 732 every new edit, and be available all the time—a fascinating  
 733 example of (almost impossible) ‘real time’ computing.  
 734 Because of all these problems a reputational system inside  
 735 Wikipedia is not near implementation as yet.

**Collective monitoring within Wikipedia: interpretation** 736

737 This Wikipedian campaign against vandalism thus has  
 738 gigantic proportions. At its core, ordinary users are called  
 739 upon to be vigilant. Meanwhile, bots have been developed  
 740 that are deployed to perform the same task. Furthermore,  
 741 users are (to some extent) provided with the tools to  
 742 implement that vigilance more effectively by combining  
 743 detection and action in one interface. Finally, the incor-  
 744 poration of a vandalism detection ‘engine’ into such tools  
 745 promises to yield the most effective anti-vandalism  
 746 approach (‘assisted editing’).  
 747

748 How to interpret this mechanism of collective  
 749 monitoring which serves to keep alive the preferred  
 750 approach of read- and write-access for all (as an expression  
 751 of full institutional trust)? A first approach is a comparison  
 752 with the notion of discretion and its associated timespan, as  
 753 developed decades ago by Elliott Jaques with a view to  
 754 determining wage levels (Jaques 1956). For (industrial)  
 755 organizations he proposed as the central characteristic of  
 756 work the amount of discretion a worker is granted, which  
 757 refers to the exercise of one’s own skills and judgment.  
 758 Discretion depends on two composing factors (Jaques  
 759 1956: Ch. III): the work content that is actually left to one’s  
 760 discretion and the procedures along which one’s perfor-  
 761 mance is reviewed by one’s superior. The longer a worker  
 762 may proceed without his boss feeling the need to review his  
 763 work (whether directly, or indirectly by observing  
 764 responses from clients), the more discretion he enjoys. This  
 765 is what Jaques refers to as the ‘time-span of discretion’  
 766 technique: determine the period of time that a worker may  
 767 enjoy without his boss intervening (ranging from weeks to  
 768 months). The emphasis here is on the *accomplishment of*  
 769 *one’s tasks*. For tasks that have a very short cycle he pro-  
 770 posed another measure of discretion (Jaques 1956: Ch. V):  
 771 the amount of ‘scrap’ a worker can produce before his boss  
 772 checks (and most probably intervenes), scrap referring to  
 773 substandard results, slower tempo, or outright damage to  
 774 tools. The more scrap from his subordinates a boss exposes  
 himself to without checking upon their performance, the

775 more discretion he may be said to grant. The emphasis in  
776 this second approach has shifted to the *avoidance of*  
777 *damage* to the resources entrusted to the worker.<sup>25</sup>

778 Jaques' apparatus can be used to illuminate the  
779 monitoring efforts within Wikipedia. Normally, users are  
780 allowed to contribute (full write-access) without any spe-  
781 cific reviewing moment on the agenda. It may take days,  
782 weeks, or months before any other user comes along and  
783 performs a check on new edits. In the wake of rising  
784 vandalism it was then decided to perform checks as soon as  
785 possible after edits were made. In terms of Jaques' second  
786 definition (damage avoidance): one came to the realization  
787 that the wikispaces had been entrusted to users without any  
788 mechanism to avoid damage being inflicted on them. Users  
789 could simply pollute, mangle, or deface Wikipedia entries  
790 considerably, in a very short time span at that. In response,  
791 the time span left to users to damage entries has drastically  
792 been reduced: from an indeterminate time span to minutes,  
793 or days at most. Users keep their discretion (full write-  
794 access), but are put under almost immediate scrutiny  
795 afterwards with every edit they perform. In that sense,  
796 users' discretion has been reduced; not by eliminating any  
797 of the tasks they are allowed to carry out, but by the  
798 introduction of much faster performance review.

799 Notice as an aside, that in this monitoring campaign  
800 ordinary users are called upon to participate without dis-  
801 tinction. They are trusted to be able and willing to carry out  
802 these tasks, besides their usual contributions. Full write-  
803 access simply includes full 'correction-access'. Ironically,  
804 this trusting gesture finds its limits as far as the supporting  
805 tools are concerned. Whenever these become stronger (and  
806 can produce more damage to Wikipedia namespaces  
807 accordingly), hesitations creep in. Strong tools such as  
808 Huggle and STiki have never been made available to the  
809 common Wikipedia, but are only granted by special per-  
810 mission (akin to rollback permission). So while the task of  
811 monitoring is unconditionally allowed to any Wikipedia, the  
812 facilitating tools for it are guarded closely and not dis-  
813 tributed indiscriminately. Trust extends to performing the  
814 *action*, not necessarily to its *instrumentation*.

815 The essential *differences* between the backgrounding of  
816 trust (by means of constant monitoring) and the substitu-  
817 tion of trust (by introducing hierarchy) can now be spelled  
818 out—note that the following argument not only applies to  
819 text/picture edits (Wikipedia) but also to source code edits  
820 (OSS). In the former approach, discretion of users is cur-  
821 tailed by the introduction of (very) frequent checks and  
822 controls; their powers of editing as such remain unchanged.

Moreover, these checks can be performed by any user with  
write-access. It is (heavy) peer-to-peer review that is  
introduced, without any hierarchical distinctions. In the  
latter approach (whether filtering/reviewing in Wikipedia,  
or the onion model in OSS), users' discretion is also  
reduced, but in a different, more drastic fashion: the grant  
of full write-access is revoked, thereby reducing the *con-*  
*tent* of their discretion. Henceforth they may only suggest  
edits, not commit them anymore. Furthermore, their edits  
are no longer scrutinized by their peers, but by a special  
layer of reviewers/committers. A division of roles has crept  
in.

A further difference between the two mechanisms  
(backgrounding vs. substituting trust) is the type of  
acceptance of text edits/code patches that is involved in  
their editorial policies. In the case of close watch,  
'*acceptance without belief*' (Cohen 1989) is and remains  
the default. That is, contributions are just accepted as-is  
and publicly displayed, without any implication that these  
can be believed. This happens for the practical purpose of  
maximizing the volume of edits coming in. Subsequent  
monitoring then may result in rejection of edits deemed  
unsuitable ('unacceptance'); otherwise, reviewed edits just  
remain accepted. In the case of hierarchy, edits are just  
considered as suggestions that may be acted upon—*non-*  
*acceptance* of edits is the default. They are simply put on  
hold (*in equilibrio*) to be reviewed. Outcomes of the sub-  
sequent review process in order to establish whether they  
can be believed or not can be twofold: either rejection or  
acceptance (as true belief). Subsequently, the edits are  
moved out of the on-hold location, and, respectively, either  
deleted, or committed to the official repository of text/code.  
The official version, therefore, always consists of *vetted*  
contributions that carry some guarantee of true belief.<sup>26</sup>

Finally, an interesting comparison can be drawn with the  
field of cyber security studies in general.<sup>27</sup> The contrasting  
approaches of backgrounding trust versus substituting trust  
can be interpreted as *access control* models. In general,  
access to websites can be controlled in several ways. On  
the one hand, lists can be drawn up that regulate access.  
With blacklisting, everyone obtains access, except those on  
a (black) list. With whitelisting the reverse obtains: the  
default is denial of access, only those on a (white) list may  
enter. On the other hand, access can be based on specific  
roles and/or permissions. The trust approaches discussed  
above can be seen to represent a *mixture* of the two logics.

25FL01 <sup>25</sup> Notice that I followed Jaques' convention of exclusive use of the  
25FL02 masculine personal pronoun throughout. Of course, it was 1956 then;  
25FL03 nowadays a gender neutral use of pronouns is considered more  
25FL04 appropriate.

26FL01 <sup>26</sup> This corresponds to the debate about models of how the human  
26FL02 mind accepts information: a procedure à la Spinoza versus a  
26FL03 procedure à la Descartes (Gilbert et al. 1990).

27FL01 <sup>27</sup> Giving ample references here would take me too far afield, but let  
27FL02 me just mention two of them. For access control, cf. O'Connor and  
27FL03 Loomis (2010); for anti-intrusion systems, cf. Scarfone and Mell  
27FL04 (2007).

869 Backgrounding trust departs from a blacklisting logic:  
 870 every user may contribute, except those on the list of  
 871 currently blocked vandals. Substituting trust departs from  
 872 the other logic, the logic of roles/permissions. By default,  
 873 edits from users are not accepted (but put ‘on hold’ for later  
 874 treatment). One category though does get unqualified  
 875 access: users who have acquired the right of auto-review  
 876 (or auto-patrol). That access is based on acquired permis-  
 877 sion, not on a (white) list of enumerated entrants.

878 Further, essential tools for the approach of backgrounding  
 879 trust are the anti-vandalism bots as mentioned above. These  
 880 in turn can be seen as an instance of an *anti-intrusion* system  
 881 for networks. Originally such systems only detected anom-  
 882 alous activities (‘intrusion detection’), leaving the corrective  
 883 action to human operators. In recent years, though, due to a  
 884 large increase in the amount of data and number of processes  
 885 involved, such systems have come to include additional  
 886 options for taking action and remedying the situation  
 887 (‘intrusion prevention’). Clearly, the Wikipedia software  
 888 bots fall into the extended category of intrusion prevention  
 889 systems, since they may—within limits—carry out actions  
 890 autonomously.<sup>28</sup>

891 **Conclusions**

892 This research has been carried out in order to answer the  
 893 following question: ‘By means of what mechanisms do co-  
 894 creative communities without restrictions on entry manage  
 895 the trusting gesture towards ‘their’ crowds of write-access  
 896 to their repositories of content?’ For easy reference,  
 897 Table 1 summarizes the results in concise terms.

898 For one thing, OSS communities were explored briefly,  
 899 warranting the conclusion that hierarchy (the ‘onion’  
 900 model) continues to apply (*substitution* of trust), the ranks  
 901 of which are filled by checking the appropriate skills of  
 902 applicants (*inference* of trust). For another, encyclopedic  
 903 communities were explored more extensively. It was found  
 904 that hierarchical roles (a form of *substituting* trust) are  
 905 created in a minimal fashion only—only one layer of  
 906 authority is usually created at the top (variously denoted  
 907 curators, constables, or administrators). In a similar vein,  
 908 the introduction of reviewing (or filtering) systems—  
 909 another instance of the *substitution* of trust—has largely  
 910 been eschewed. Instead, the requirement of registration in  
 911 some encyclopedias serves the purpose of being able to  
 912 *infer* some degree of trustworthiness on the part of users.

28FL01 <sup>28</sup> Another comparison of an epistemological kind can be drawn. The  
 28FL02 Wikipedia monitoring campaign can be seen as an institutional form  
 28FL03 of ‘epistemic vigilance’ concerning information communicated by  
 28FL04 others—as Sperber et al. (2010) coined the term. In our case, such  
 28FL05 vigilance is not exercised in judicial or scientific institutions (idem:  
 28FL06 383), but in the largest open-content encyclopedia of all.

913 Wikipedia however, the largest encyclopedia of all, is  
 914 plagued by persistent vandalism. As a line of defence  
 915 against vandals, a *backgrounding* mechanism is resorted to:  
 916 collective monitoring. Forms and intensity of this constant  
 917 ‘peer review’ of fresh edits within Wikipedia have grown  
 918 considerably. In addition, software bots have been called to  
 919 the rescue. By now the task of monitoring is divided  
 920 equally between humans and bots. Its essence is keeping  
 921 the discretionary powers of contributing Wikipedians  
 922 unchanged, while focusing upon damage avoidance by  
 923 shortening the times of review as executed by the com-  
 924 munity as a whole. Not discretionary content itself, but the  
 925 associated damage potential is curtailed.

926 Further, the envisaged role of reputation in such ‘epi-  
 927 stemic vigilance’ has been explored (the term is borrowed  
 928 from Sperber et al. 2010; cf. note 28). As far as suitable  
 929 indicators of reputation can be constructed, these can  
 930 generally be used for distributing privileges or assigning  
 931 roles related to ‘production’. By a similar reasoning, ideas  
 932 circulate in Wikipedia of incorporating (a measure of)  
 933 reputation—as indicator from which the trustworthiness of  
 934 contributors can be *inferred*—into the very heart of col-  
 935 lective monitoring itself. The quality of fresh edits is to be  
 936 gauged by the momentary reputational score of their con-  
 937 tributors, thus presumably enhancing the efficiency of  
 938 vandalism detection. Just as hierarchy is usually inter-  
 939 twined with proven track record, vandalism detection can  
 940 be entwined with reputation in order to achieve optimal  
 941 outcomes.

942 It is worth remarking that the ‘backgrounded’ model  
 943 would seem to represent the most elaborated efforts with  
 944 which an institutional policy of ‘acceptance without belief’  
 945 can be kept afloat. If the campaign of constantly monitor-  
 946 ing content cannot back up institutional trust firmly  
 947 enough, the Wikipedia repository of knowledge may  
 948 gradually get corrupted: its contents become polluted and  
 949 its integrity is compromised. If so, steps towards bureau-  
 950 cracy (substituting trust) seem to be the only option left.  
 951 Wikipedia governance would then inescapably move in  
 952 the direction of the type of governance customary for OSS.

953 As clearly can be seen now, OSS and Wikipedia  
 954 models of governance differ essentially in the way in which  
 955 institutional trust is granted. In the former model, trust is  
 956 only granted *ex post*: to the extent that one has proven to be  
 957 loyal and capable, one may obtain ever more permissions  
 958 and rise to higher levels in the hierarchy. In the Wikipedia  
 959 model, however, considerable trust is already granted *ex*  
 960 *ante*: full write-access is (almost immediately) open to  
 961 anyone. In close connection, incessant ‘peer review’ is  
 962 exercised to keep this option viable and feasible.

963 Finally, a speculation is in order. The mode of ‘epistemic  
 964 vigilance’ fuelled by reputation tracking may be general-  
 965 ized into a model of how the institution of Wikipedia *as a*

**Table 1** Mechanisms used to handle the institutional assumption of trust in contributors (fully open co-creative open-content communities)

	OSS	Scholarpedia <sup>a</sup>	Citizendium	Wikinews	Wikipedia
Inferring trust	Track record	Registration	Registration		Reputation <sup>b</sup>
Backgrounding trust					Collective monitoring
Substituting trust	Role hierarchy (onion model)	'Curators'	'Constables'	'Administrators'	'Administrators' (provided with supplementary 'permissions')
				New edit review <sup>c</sup>	New edit review <sup>d</sup>

<sup>a</sup> Only applies to the publication phase

<sup>b</sup> Under consideration

<sup>c</sup> Only used for the 'front page' ('flagged-revisions' scheme)

<sup>d</sup> Only used for entries subjected to persistent vandalism ('pending-changes' protection)

966 whole might operate in the future (cf. also de Laat 2012a).  
 967 Reputations of contributors are constantly monitored and  
 968 updated. In a 'big data' process, each and every vote  
 969 increases or decreases the credibility of entries and updates  
 970 reputations of contributors on a continuous basis. The  
 971 measure of reputation is no longer the gauge for collective  
 972 monitoring only, but for all governance: for the distribution  
 973 of privileges (such as permissions to use high-powered  
 974 'assisted editing' tools) and the membership of hierarchical  
 975 layers (such as 'reviewers'). Consequently, one's fate is  
 976 constantly in the balance. When one's reputation, as  
 977 indicator of trustworthiness, rises, one may obtain more  
 978 privileges and rise in the ranks. But a declining reputation  
 979 (signalling a decrease of trustworthiness) necessitates a  
 980 close watch on one's activities; continued misbehaviour  
 981 may lead to a loss of privilege and/or role access, and  
 982 ultimately result in expulsion from the community. As can  
 983 be seen, in the model a sizeable basic trust level is the  
 984 starting point (*ex ante*), to be gauged and tested continu-  
 985 ously as time develops.

986 In closing, let me remark that the above analysis can  
 987 usefully be regarded as the necessary groundwork for a  
 988 'disclosive' ethical approach (in the sense of Brey 2000)  
 989 towards Wikipedia in particular. Hidden characteristics of  
 990 its governance have been laid bare that may have moral  
 991 connotations. One may proceed now to a more fully  
 992 informed ethical analysis of the encyclopedia's governance,  
 993 focusing on the uncovered mechanism of backgrounding  
 994 trust. Essential values like privacy, justice, and even  
 995 democracy seem to be implicated. What are the moral  
 996 merits of close surveillance zooming in on anonymous  
 997 contributors? Is it morally justified to let supposed 'vandals'  
 998 be repelled by bots that proceed autonomously? What levels  
 999 of false positives can justifiably be chosen for the actions of  
 1000 such bots? How can Wikipedia justify to its visitors the  
 1001 tension between open write-access on the one hand and  
 1002 constant monitoring on the other? Does the practice of  
 1003 close surveillance by any chance amount to abandoning  
 1004 Wikipedia's original mission of a 'democratic' production

of knowledge? Such questions need to be answered urgently 1005  
 (for some answers, cf. de Laat 2014). 1006

**Acknowledgments** Thanks are due to two anonymous reviewers of 1007  
 this journal for their comments, in particular for alerting me to the 1008  
 comparison with cyber security studies. 1009

## References 1010

*All websites were last accessed on February 10, 2014.* 1011

- Adler, B. T., Chatterjee, K., de Alfaro, L., Faella, M., Pye, I., & Raman, V. (2008). Assigning trust to Wikipedia content. In *Proceedings of the 4th International Symposium on Wikis (WikiSym'08)*, September 8–10, 2008, Porto, Portugal. <http://dx.doi.org/10.1145/1822258.1822293>. 1012  
 1013  
 1014  
 1015  
 1016  
 Adler, B. T., & de Alfaro, L. (2007). A content-driven reputation system for the Wikipedia. In *Proceedings of the 16th International Conference on World Wide Web*, May 8–12, 2007, Banff, Alberta, Canada. <http://dx.doi.org/10.1145/1242572.1242608>. 1017  
 1018  
 1019  
 1020  
 Adler, B. T., de Alfaro, L., Mola-Velasco, S. M., Rosso, P., & West, A. G. (2011). Wikipedia vandalism detection: Combining natural language, metadata, and reputation features. In *CICLing '11: Proceedings of the 12th International Conference on Intelligent Text Processing and Computational Linguistics, LNCS 6609* (pp. 277–288), Tokyo, Japan. 1021  
 1022  
 1023  
 1024  
 1025  
 1026  
 Brey, P. (2000). Disclosive computer ethics. *Computers and Society*, 30(4), 10–16. 1027  
 1028  
 Cohen, L. J. (1989). Belief and acceptance. *Mind, New Series*, 98(391), 367–389. 1029  
 1030  
 Crowston, K., Annabi, H., Howison, J., & Masango, Ch. (2004). Effective work practices for software engineering: Free/libre open source software development. In *Proceedings of the 2004 ACM workshop on Interdisciplinary software engineering research (WISER '04)* (pp. 18–26), ACM, New York, NY, USA. <http://doi.acm.org/10.1145/1029997.1030003>. 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 de Laat, P. B. (2007). Governance of open source software: State of the art. *Journal of Management and Governance*, 11(2), 165–177. 1037  
 1038  
 1039  
 de Laat, P. B. (2010). How can contributors to open-source communities be trusted? On the assumption, inference, and substitution of trust. *Ethics and Information Technology*, 12(4), 327–341. 1040  
 1041  
 1042  
 1043



- 1044 de Laat, P. B. (2012a). Open source production of encyclopedias:  
1045 Editorial policies at the intersection of organizational and  
1046 epistemological trust. *Social Epistemology*, 26(1), 71–103.
- 1047 de Laat, P. B. (2012b). Coercion or empowerment? Moderation of  
1048 content in Wikipedia as ‘essentially contested’ bureaucratic  
1049 rules. *Ethics and Information Technology*, 14(2), 123–135.
- 1050 de Laat, P. B. (2012c). Navigating between chaos and bureaucracy:  
1051 Backgrounding trust in open-content communities. In K. Aberer  
1052 et al. (Eds.), *Proceedings of the 4th International Conference on*  
1053 *Social Informatics, SocInfo 2012, LNCS 7710, Heidelberg:*  
1054 *Springer* (pp. 534–557), December 5–7, Lausanne, Switzerland.
- 1055 de Laat, P. B. (2014). Tools and bots against vandalism: Eroding  
1056 Wikipedia’s moral order? In *Proceedings of the 11th Interna-*  
1057 *tional Conference of Computer Ethics: Philosophical Explora-*  
1058 *tions (CEPE)*, Paris.
- 1059 Deterding, S., Dixon, D., Khaled, R., & Nacke, L. E. (2011). From  
1060 game design elements to gamefulness: Defining »Gamification«.  
1061 In *Mindtrek 2011 Proceedings*, Tampere: ACM Press.
- 1062 Dutton, W. H. (2008). The wisdom of collaborative network  
1063 organizations: Capturing the value of networked individuals.  
1064 *Prometheus*, 26(3), 211–230.
- 1065 Farmer, F. R., & Glass, B. (2010). *Building web reputation systems*.  
1066 Sebastopol: O’Reilly.
- 1067 Forte, A., & Lampe, C. (2013). Defining, understanding, and  
1068 supporting open collaboration: Lessons from the literature.  
1069 *American Behavioral Scientist*, 57(5), 535–547.
- 1070 Gilbert, D., Krull, D., & Malone, P. (1990). Unbelieving the  
1071 unbelievable: Some problems in the rejection of false information.  
1072 *Journal of Personality and Social Psychology*, 59(4), 601–613.
- 1073 Holck, J., & Jørgensen, N. (2005). Do not check in on red: Control  
1074 meets anarchy in two open source projects. In S. Koch (Ed.),  
1075 *Free/open source software development* (pp. 1–26). Hershey:  
1076 Idea Group.
- 1077 Jaques, E. (1956). *Measurement of responsibility: A study of work,*  
1078 *payment, and individual capacity*. London: Tavistock.
- 1079 McGeer, V. (2008). Trust, hope and empowerment. *Australasian*  
1080 *Journal of Philosophy*, 86(2), 237–254.
- Moran, R. (2005). Getting told and being believed. *Philosophers’*  
1081 *Imprint*, 5(5); also published in J. Lackey, & E. Sosa (Eds.)  
1082 (2006), *The epistemology of testimony* (pp. 272–306). Oxford:  
1083 Oxford University Press.
- O’Connor, A. C., & Loomis, R. J. (2010). *Economic analysis of role-*  
1084 *based access control*. Prepared for NIST. [http://csrc.nist.gov/](http://csrc.nist.gov/groups/SNS/rbac/documents/20101219_RBAC2_Final_Report.pdf)  
1085 [groups/SNS/rbac/documents/20101219\\_RBAC2\\_Final\\_Report.](http://csrc.nist.gov/groups/SNS/rbac/documents/20101219_RBAC2_Final_Report.pdf)  
1086 [pdf](http://csrc.nist.gov/groups/SNS/rbac/documents/20101219_RBAC2_Final_Report.pdf).
- Orsila, H., Geldenhuys, J., Ruokonen, A., & Hammouda, I. (2009).  
1087 Trust issues in open source software development. In N.  
1088 Medvidovic, & T. Tamai (Eds.), *Proceedings of the Warm Up*  
1089 *Workshop for ACM/IEEE ICSE 2010 (WUP ’09)* (pp. 9–12),  
1090 ACM, New York, NY, USA. [http://doi.acm.org/10.1145/](http://doi.acm.org/10.1145/1527033.1527037)  
1091 [1527033.1527037](http://doi.acm.org/10.1145/1527033.1527037).
- Pullinger, K. (2012). A million penguins’ five years on, blog post  
1092 from 25 January 2012. [http://www.katepullinger.com/blog/com-](http://www.katepullinger.com/blog/comments/a-million-penguins-five-years-on/)  
1093 [ments/a-million-penguins-five-years-on/](http://www.katepullinger.com/blog/comments/a-million-penguins-five-years-on/).
- Scarfone, K. and Mell, P. (2007). *Guide to intrusion detection and*  
1094 *prevention systems (IDPS)*. NIST Special Publication 800-94.  
1095 <http://csrc.nist.gov/publications/nistpubs/800-94/SP800-94.pdf>.
- Sperber, D., Clément, F., Heintz, C., Mascaro, O., Mercier, H.,  
1096 Origgi, G., & Wilson, D. (2010). Epistemic vigilance. *Mind &*  
1097 *Language*, 25, 359–393.
- West, A. G. (2011). Anti-vandalism research: The year in review  
1098 (Presentation at Wikimania 2011). [www.cis.upenn.edu/](http://www.cis.upenn.edu/~westand/docs/wikimania_11_vandalism_slides.pdf)  
1099 [~westand/docs/wikimania\\_11\\_vandalism\\_slides.pdf](http://www.cis.upenn.edu/~westand/docs/wikimania_11_vandalism_slides.pdf).
- West, A. G., Chang, J., Venkatasubramanian, K. K., & Lee, I. (2012).  
1100 Trust in collaborative web applications. *Future Generation*  
1101 *Computer Systems*, 28(8), 1238–1251. [http://dx.doi.org/10.1016/](http://dx.doi.org/10.1016/j.future.2011.02.007)  
1102 [j.future.2011.02.007](http://dx.doi.org/10.1016/j.future.2011.02.007).
- WP:Huggle. <http://en.wikipedia.org/wiki/Wikipedia:Huggle>.  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117