# Fixing Dockerfile Smells: An Empirical Study

Giovanni Rosa  $\cdot$  Federico Zappone  $\cdot$  Simone Scalabrino  $\cdot$  Rocco Oliveto

Received: date / Accepted: date

Abstract Docker is the *de facto* standard for software containerization. A 1 Dockerfile contains the requirements to build a Docker image containing a tar-2 get application. There are several best practice rules for writing Dockerfiles, 3 but the developers do not always follow them. Violations of such practices, 4 known as Dockerfile smells, can negatively impact the reliability and perfor-5 mance of Docker images. Previous studies showed that Dockerfile smells are 6 widely diffused, and there is a lack of automatic tools that support developers 7 in fixing them. However, it is still unclear what Dockerfile smells get fixed by 8 developers and to what extent developers would be willing to fix smells in the 9 first place. The aim of our study is twofold. First, we want to understand what 10 Dockerfiles smells receive more attention from developers, *i.e.*, are fixed more 11 frequently in the history of open-source projects. Second, we want to check if 12 developers are willing to accept changes aimed at fixing Dockerfile smells (e.q., 13 generated by an automated tool), to understand if they care about them. We 14 evaluated the survivability of Dockerfile smells from a total of 53,456 unique 15 Dockerfiles, where we manually validated a large sample of smell-removing 16 commits to understand (i) if developers performed the change with the inten-17 tion of removing bad practices, and (ii) if they were aware of the removed smell. 18

G. Rosa STAKE Lab, University of Molise, Italy E-mail: giovanni.rosa@unimol.it

F. Zappone STAKE Lab, University of Molise, Italy E-mail: f.zappone1@studenti.unimol.it

S. Scalabrino

STAKE Lab, University of Molise, Italy E-mail: simone.scalabrino@unimol.it

R. Oliveto

STAKE Lab, University of Molise, Italy E-mail: rocco.oliveto@unimol.it In the second part, we used a rule-based tool to automatically fix Dockerfile smells. Then, we proposed such fixes to developers via pull requests. Finally, we quantitatively and qualitatively evaluated the outcome after a monitoring period of more than 7 months. The results of our study showed that most developers pay more attention to changes aimed at improving the performance of Dockerfiles (image size and build time). Moreover, they are willing to accept the fixes for the most common smells, with some exceptions (*e.g.*, missing version pinning for OS packages).

Keywords dockerfile smells · empirical software engineering · software
 evolution

## 11 1 Introduction

Software systems are developed to be deployed and used. Operating software 12 in a production environment, however, entails several challenges. Among the 13 others, it is very important to make sure that the software system behaves 14 exactly as in a development environment. Virtualization and, above all, con-15 tainerization technologies are increasingly being used to ensure that such a 16 requirement is  $met^1$ . Among the others,  $Docker^2$  is one of the most popular 17 platforms used in the DevOps workflow: It is the main containerization frame-18 work in the open-source community [6], and is widely used by professional 19 developers<sup>3</sup>. Also, Docker is the most loved and most wanted platform in 20 the 2021 StackOverflow survey<sup>3</sup>. Docker allows releasing applications together 21 with their dependencies through containers (*i.e.*, virtual environments) shar-22 ing the host operating system kernel. Each Docker image is defined through a 23 Dockerfile, which contains instructions to build the image containing the appli-24 cation. All the public Docker images are hosted on an online repository called 25 DockerHub<sup>4</sup>. Since its introduction in 2013, Docker counts 3.3M of Desktop 26 installations, and 318B image pulls from DockerHub<sup>5</sup>. 27 Defining Dockerfiles, however, is far from trivial: Each application has its 28 own dependencies and requires specific configurations for the execution envi-29 ronment. Previous work [21] introduced the concept of Dockerfile smells, which 30

<sup>30</sup> are violations of best practices, similarly to code smells [5], and a catalog of

 $_{32}$  such problems<sup>6</sup>. The presence of such smells might increase the risk of build

failures, generate oversized images, and security issues [6, 10, 22, 23]. Previous

work studied the prevalence of Dockerfile smells [6,9,14].

Despite the popularity and adoption of Docker, there is still a lack of tools to support developers in improving the quality and reliability of containerized

 $<sup>^1</sup>$  https://portworx.com/blog/2017-container-adoption-survey/

 $<sup>^2</sup>$  https://www.docker.com/

 $<sup>^3</sup>$  https://insights.stackoverflow.com/survey/2021

<sup>4</sup> https://hub.docker.com/

<sup>&</sup>lt;sup>5</sup> https://www.docker.com/company/

<sup>&</sup>lt;sup>6</sup> https://github.com/hadolint/hadolint/wiki

<sup>1</sup> applications, e.g., tools for automatic refactoring of code smells on Docker-

<sup>2</sup> files [13]. Relevant studies in this area investigated the prevalence of Dockerfile

 $_{3}$  smells in open-source projects [6,9,14,21], the diffusion technical debt [4], and

<sup>4</sup> the refactoring operations typically performed by developers [13].

<sup>5</sup> While it is clear which Dockerfile smells are more frequent than others, it is <sup>6</sup> still unclear which smells are more important to developers. A previous study <sup>7</sup> by Eng *et al.* [9] reported how the number of smells evolves over time. Still, <sup>8</sup> there is no clear evidence showing that (i) developers actually fix Dockerfile <sup>9</sup> smells (*e.g.*, they might incidentally disappear), and that (ii) developers would <sup>10</sup> be willing to fix Dockerfile smells in the first place.

In this paper, we propose a study to fill this gap. First, we analyze the 11 survivability of Dockerfile smells to understand how developers fix them and 12 which smells they consider relevant to remove. This, however, only tells a part 13 of the story: Developers might not correct some smells because they are harder 14 to fix. Therefore, we also evaluated to what extent developers are willing to 15 accept fixes to smells when they are proposed to them (e.g., by an automated16 tool). The context of the study is represented by a total of 220k commits 17 and 4,255 repositories, extracted from a state-of-the-art dataset containing 18 the change history of about 9.4M unique Dockerfiles. 19

For each instance of such a dataset (which is a Dockerfile snapshot), we 20 extracted the list of Dockerfile smells using the *hadolint* tool [2]. The tool 21 performs a rule check on a parsed Abstract Syntax Tree (AST) representation 22 of the input Dockerfile, based on the Docker [1] and shell script [3] best prac-23 tices. Next, we manually validate a total of 1,000 commits that make one or 24 more smells disappear to verify (i) that they are real fixes (e.q.), the smell was 25 not removed incidentally), (ii) whether the fix is *informed* (e.g., if developers 26 explicitly mention such an operation in the commit message), and (iii) remove 27 possible false positives identified by hadolint. 28

Then, we evaluated to what extent developers are willing to accept changes 29 aimed at fixing smells. To this aim, we defined DOCKLEANER, a rule-based 30 refactoring tool that automatically fixes the 12 most frequent Dockerfile smells. 31 We used DOCKLEANER to fix a set of smelly Dockerfiles extracted from the 32 most active repositories. Next, we submitted a total of 157 pull requests to de-33 velopers containing the fixes, one for each repository. We monitored the status 34 of the pull requests for more than 7 months (*i.e.*, 218 days). In the end, we eval-35 uated how many of them get accepted for each smell type and the developers' 36 reactions. The results show that, mostly, smells are fixed either very shortly 37 (36% of the cases). There are also cases in which they are fixed after a very 38 long period (2% - after 2 years). This could be a consequence of the fact that, 39 generally, a few changes are performed on Dockerfiles and there the probabil-40 ity of noticing the errors is higher in the short-term (e.g., until the Dockerfile 41 works correctly) or, instead, it naturally increases with time, but very slowly. 42 Also, developers perform changes on Dockerfiles mainly to optimize the build 43 time and reduce the final image size, while there are only few changes limited 44 only to the improvement of code quality. Even if Dockerfile smells are com-45 monly diffused among Dockerfiles, developers are gradually becoming aware 46

of the writing best practices for Dockerfiles. For example, avoiding the usage 1 of MAINTAINER which is deprecated, or they prefer to use COPY instead of ADD 2 for copying files and folders as it is suggested by the Docker guidelines<sup>7</sup>. In 3 addition, developers are open to approve changes aimed at fixing smells for the 4 most common violations, but with some exceptions. Examples are the missing 5 version pinning for apt-get packages (DL3008), which has received negative 6 reactions from developers. However, version pinning, in general, is considered 7 fundamental for other aspects, such as the base image pinning (DL3006 and 8 DL3007), or the pinning of software dependencies (e.g., npm and pip). 9 To summarize, the contributions that we provided with our study are the 10

- 11 following:
- We performed a detailed analysis of the survivability of Dockerfile smells
   and manually validated a sample of smell-fixing commits for Dockerfile
   smells;
- We introduced DOCKLEANER, a rule-based tool to fix the most common
   Dockerfile smells;
- 3. We ran an evaluation via pull requests of the willingness of developers of
   accepting changes aimed at fixing Dockerfile smells.

The remaining of the paper is organized as follows: In Section 2 we provide a general overview on Dockerfile smells and related works. Section 3 describes the

design of our study, while in Section 5 we present the results of our experiment.

<sup>22</sup> In section Section 6 we qualitatively discuss the results. Finally, Section 7

<sup>23</sup> discusses the threats to validity and in Section 8 we summarize some final

24 remarks and future directions.

## 25 2 Background and Related Work

Technical debt [12] has a negative impact on the software maintainability. A 26 symptom of technical debt is represented by *code smells* [5]. Code smells are 27 poor implementation choices, that does not follow design and coding best prac-28 tices, such as design patterns. They can negatively impact the maintainability 29 of the overall software system. Mainly, code smells are defined for object-30 oriented systems. Some examples are duplicated code or god class (*i.e.*, a class 31 having too much responsibilities). In the following, we first introduce smells 32 that affect Dockerfile, and then we report recent studies on their diffusion and 33 the practices used to improve Dockerfile quality. 34

Dockerfile smells. Docker reports an official list of best practices for writing Dockerfiles [1]. Such best practices also include indications for writing shell
 script code included in the RUN instructions of Dockerfiles. For example, the
 usage of the instruction WORKDIR instead of the bash command cd to change

<sup>39</sup> directory. This because each Docker instruction defines a new layer at the time

 $<sup>^7</sup>$  https://docs.docker.com/develop/develop-images/dockerfile\_best-practices/\# add-or-copy

of build. The violation of such practices lead to the introduction of Dockerfile 1 smells. In fact, with Dockerfile smells, we indicate that instructions of a Dock-2 erfile that violate the writing best practices and thus can negatively affect the 3 quality of them [21]. The presence of Dockerfile smells can also have a direct 4 impact on the behavior of the software in a production environment. For exam-5 ple, previous work showed that missing adherence to best practices can lead to 6 security issues [22], negatively impact the image size [10], increase build time 7 and affect the reproducibility of the final image (*i.e.*, build failures) [6, 10, 23]. 8 For example, the version pinning smell, that consists in missing version num-9 ber for software dependencies, can lead to build failures as with dependencies 10 updates the execution environment can change. There are several tools that 11 support developers in writing Dockerfiles. An example is the *binnacle* tool, pro-12 posed by Henkel et al. [10] that performs best practices rule checking defined 13 on the basis of a dataset of Dockerfiles written by experts. The reference tool 14 used in literature for the detection of Dockerfile smells is hadolint [2]. Such 15 a tool checks a set of best practices violations on a parsed AST version of 16 the target Dockerfile using a rule-based approach. Hadolint detects two main 17 categories of issues: Docker-related and shell-script-related. The former affect 18 Dockerfile-specific instructions (e.g., the usage of absolute path in the WORKDIR 19 command<sup>8</sup>). They are identified by a name having the prefix DL followed by 20 a number. The shell-script-related violations, instead, specifically regard the 21 shell code in the Dockerfile (e.q., in the RUN instructions). Such violations are 22 a subset of the ones detected by the ShellCheck tool [3] and they are identified 23 by the prefix SC followed by a number. It is worth saying that these rules 24 can be updated and changed during time. For example, as the instruction 25 MAINTAINER has been deprecated, the rule DL4000 that previously check for 26 the usage of that instructions that was a best practice, has been updated as 27 the avoidance of that instruction because it is deprecated. 28

Diffusion of Dockerfile smells. A general overview of the diffusion of 29 Dockerfile smells was proposed by Wu et al. [21]. They performed an empirical 30 study on a large dataset of 6.334 projects to evaluate which Dockerfile smells 31 occurred more frequently, along with coverage, distribution and a particular 32 focus on the relation with the characteristics of the project repository. They 33 found that nearly 84% of GitHub projects containing Dockerfiles are affected 34 by Dockerfile smells, where the Docker-related smells are more frequent that 35 the shell-script smells. Also in this direction, Cito et al. [6] performed an em-36 pirical study to characterize the Docker ecosystem in terms of quality issues 37 and evolution of Dockerfiles. They found that the most frequent smell regards 38 the lack of version pinning for dependencies, that can lead to build fails. Lin 39 et al. [14] conducted an empirical analysis of Docker images from Docker-40 Hub and the git repositories containing their source code. They investigated 41 different characteristics such as base images, popular languages, image tag-42 ging practices and evolutionary trends. The most interesting results are those 43 related to Dockerfile smells prevalence over time, where the version pinning 44

<sup>&</sup>lt;sup>8</sup> https://github.com/hadolint/hadolint/wiki/DL3000

smell is still the most frequent. On the other hand, smells identified as DL3020 1 (*i.e.*, COPY/ADD usage), DL3009 (*i.e.*, clean apt cache) and DL3006 (*i.e.*, image 2 version pinning) are no longer as prevalent as before. Furthermore, violations 3 DL4006 (*i.e.*, usage of RUN pipefail) and DL3003 (*i.e.*, usage of WORKDIR) be-4 came more prevalent. Eng et al. [9] conducted an empirical study on the largest 5 dataset of Dockerfiles, spanning from 2013 to 2020 and having over 9.4 million 6 unique instances. They performed an historical analysis on the evolution of 7 Dockerfiles, reproducing the results of previous studies on their dataset. Also 8 in this case, the authors found that smells related to version pinning (i.e.,9 DL3006, DL3008, DL3013 and DL3016) are the most prevalent. In terms of 10 Dockerfile smell evolution, they show that the count of code smells is slightly 11 decreasing over time, thus hinting at the fact that developers might be inter-12 ested in fixing them. Still, it is unclear the reason behind their disappearance, 13

e.g., if developers actually fix them or if they get removed incidentally.

#### 15 3 Study Design

<sup>16</sup> The *goal* of our study is to understand whether developers are interested in

17 fixing Dockerfile smells. The *perspective* is of researchers interested in improv-

ing Dockerfile quality. The *context* consists in 53,456 Dockerfile snapshots,
 extracted from 4,255 repositories.

<sup>20</sup> In detail, the study aims to address the following research questions:

- RQ<sub>1</sub>: How do developers fix Dockerfile smells? We want to conduct a
 comprehensive analysis of the survivability of Dockerfile smells. Thus, we
 investigate what smells are fixed by developers and how.

 $_{24}$  – **RQ<sub>2</sub>:** Which Dockerfile smells are developers willing to address? We want

 $_{25}$  to understand if developers would find beneficial changes aimed at fixing

 $_{26}$  Dockerfile smells (*e.g.*, generated by an automated refactoring tool).

#### 27 3.1 Study Context

<sup>28</sup> The context of our study is represented by a subset of the dataset introduced by

<sup>29</sup> Eng et al. [9]. The dataset consists in about 9.4 million Dockerfiles, in a period

 $_{30}$   $\,$  spanning from 2013 to 2020. To the best of our knowledge, the dataset is the

 $_{31}$  largest and the most recent one from those available in the literature [6,10,13].

<sup>32</sup> Moreover, such a dataset contains the change history (*i.e.*, commits) of each

<sup>33</sup> Dockerfile. This characteristic allows us to evaluate the survivability of code <sup>34</sup> smells (RQ<sub>1</sub>). The authors constructed that dataset through mining software

<sup>35</sup> repositories from the S version of the WoC (World of Code) dataset [15].

#### 36 3.2 Data Collection

- <sup>1</sup> To avoid toy projects, we selected only the repositories having at least 10 stars
- <sup>2</sup> for a total of 4,255 repos, excluding forks. We also discarded the repositories



Fig. 1: Overall workflow of the experimentation procedure.

where the star number is not available in the original dataset (*i.e.*, the value 3 is reported as NULL). We cloned all the available repositories from the selected 4 sample to obtain the most updated commit data at the time our analysis 5 started (*i.e.*, March 2023). Next, using a heuristic approach, we (i) identified all 6 the Dockerfiles at the latest commit, and (ii) we traversed the commit history 7 to get all the commits and snapshots for the identified Dockerfile. In detail, for 8 the first step, we processed all the source files contained in the repository and 9 we evaluated if the file (i) contains the word "dockerfile" in the filename, and 10 (ii) if contains valid and non-empty commands, *i.e.*, can be correctly parsed 11 using the official  $dockerfile \ parser^9$ . For each valid Dockerfile, we mined the 12 change history using git log. We excluded the Dockerfiles having only one 13 snapshot (*i.e.*, no changes, referenced by only one commit). After this, we 14

<sup>15</sup> extracted a total of 220k commits corresponding to 53,456 unique Dockerfiles.

In the end, we ran the latest version of  $hadolint^{10}$  for each Dockerfile to extract

<sup>17</sup> the Dockerfile smells, if present.

## 18 4 Experimental Procedure

<sup>19</sup> In this section, we describe the experimentation procedure that we will use to <sup>20</sup> answer our RQs. Fig. 1 describes the overall workflow of the study.

#### 21 4.1 RQ<sub>1</sub>: How do developers fix Dockerfile smells?

<sup>1</sup> To answer RQ<sub>1</sub>, we perform an empirical analysis on Dockerfile smell surviv-

 $_{2}$  ability. For each Dockerfile d, associated with the respective repository from

<sup>3</sup> GitHub, we consider its snapshots over time,  $d_1, \ldots, d_n$ , associated with the

 $^{10}$  hadolint release v2.12.0

<sup>&</sup>lt;sup>9</sup> https://github.com/asottile/dockerfile

1	RUN apt-get install —y \	1	RUN apt-get install -y \
2	curl=7.* \	2	curl=7.* \
3	wget \		
4	&& rm -rf /var/lib/apt/lists/*	3	&& rm -rf /var/lib/apt/lists/*

Fig. 2: Example of a candidate smell-fixing commit that does not actually fix the smell.

respective commit IDs in which they were introduced  $(i.e., c(d_1), \ldots, c(d_n))$ . 4 We also consider the Dockerfile smells detected with *hadolint*, indicated as 5  $\eta(d_1), \ldots, \eta(d_n)$ . For each snapshot  $d_i$  (with i > 1) of each Dockerfile d, we 6 compute the disappeared smells as  $\delta(d_i) = \eta(d_i) - \eta(d_{i-1})$ . All the snapshots 7 for which  $\delta(d_i)$  is not an empty set are *candidate* changes that aim at fixing the 8 smells. We define a set of all such snapshot as  $PF = \{d_i : |\delta(d_i)| > 0\}$ . In the 9 end, we obtain a set of smelly  $(d_{i-1})$  and smell-removing commit  $(d_i)$  pairs. We 10 implemented the described procedure as a basic heuristic approach, which (i) 11 went through all the commits, (ii) executed *hadolint* to detect smells, (iii) re-12 turned the smelly and smell-removing commits pairs. The total time required 13 was about nine hours. 14

Next, we manually evaluate the commit pairs to verify (i) that the changes 15 that led to the snapshots in PF are actual fixes for the Dockerfile smell, and 16 (ii) whether developers were aware of the smell when they made the change, 17 and (iii) avoid any bias related to the presence of false positives in terms of 18 smells (identified by hadolint). In detail, we manually inspect a sample of 1,000 19 of such candidate changes, which is statistically representative, leading to a 20 margin of error of 3.1% (95% confidence interval) assuming an infinitely large 21 population. We look at the code diff to understand how the change was made 22 (*i.e.*, if it fixed the smell or if the smell disappeared incidentally). Also, for ac-23 tual fixes, we consider the commit message, the possible issues referenced in it, 24 and the pull requests to which they possibly belong to understand the purpose 25 of the change (*i.e.*, if the fix was informed or not). We identify as *smell fixing* 26 change a commit in which developers (i) modified one or more Dockerfile lines 27 that contained one or more smells in the previous snapshot (*i.e.*, commit), and 28 (ii) kept the functionality expressed in those lines. For example, if the commit 29 removes the instruction line where the smell is present, we do not label it as 30 an actual smell-fixing commit. This is because the smelly line is just removed 31 and not fixed (*i.e.*, the functionality changed). Let us consider the example 32 in Fig. 2: The package wget lacks version pinning (left). An actual fix would 33 consist of the addition of a version to the package. Instead, in the commit, the 34 package gets simply removed (e.g., because it is not necessary). Therefore, we 35 do not consider such a change as a fixing change. Besides, we mark a fix as 36 informed if the commit message, the possibly related pull request, or the issue 37 possibly fixed with the commit explicitly reports that the modification aimed 1

<sup>2</sup> to fix a bad practice.

8

Table 1: The most frequent Dockerfile smells identified in literature [9], along with the most fixed rules we identified in our study (reported with \*). We implemented all of the rules in DOCKLEANER.

Rule	Description	How to fix
DL3003	Use WORKDIR to switch to a directory	Replace cd command with WORKDIR
DL3006	Missing version pinning for base image	Pin the version tag corresponding to the resulting im- age digest
DL3008	Missing version pinning of <b>apt-get</b> packages	Pin the latest suitable package version from Launch- pad
DL3009	Delete the ${\tt apt-get}$ lists after installing packages	Add in the corresponding instruction block the lines to clean apt cache
DL3015	Avoid additional packages by specifyingno-install-recommends	Add the optionno-install-recommends to the cor- responding instruction block
DL3020	Use $\texttt{COPY}$ instead of $\texttt{ADD}$ for files and folders	Replace ADD instruction with COPY when copying files and folders
DL4000	MAINTAINER is deprecated	Replace maintainer with the equivalent LABEL instruc- tion
DL4006	Set -o pipefall to avoid silencing errors in RUN in- structions having pipe operations	Add the SHELL pipefail instruction before RUN that uses pipe
$DL3059^*$	Consider consolidation for multiple consecutive RUN in- structions	Concatenate all subsequent RUN instruction until a comment line or a different instruction
$DL3007^*$	Avoid to use the latest to tag the version of an image	Same approach as DL3006
$DL3025^*$	Use arguments JSON notation for CMD and ENTRYPOINT	Refactor the instruction command as JSON notation
$DL3048^*$	Invalid Label Key	Refactor the LABEL instructions according to the hadolint documentation examples <sup>11</sup>

Two of the authors independently evaluated each instance. The evalua-3 tors discussed conflicts for both the aspects evaluated aiming at reaching a 4 consensus. The agreement between the two annotators is measured using the 5 Cohen's Kappa Coefficient [7], obtaining a value of k = 0.79 considered "very" 6 good" according to the interpretation recommendations [16]. The total effort 7 required for the manual validation was about five working days, considering 8 two of the authors that performed the annotation and discussed the conflicts. 9 Moreover, starting from the smell-fixing change, we go back through the 10 change history to identify the *last-smell-introducing* commit, *i.e.*, the commit 11 in which the artifact can be considered smelly [19], by executing git blame 12 on the Dockerfile line number labeled as smelly by hadolint. In the end, we 13 summarize the total number of fix commits and the percentage of actual fix 14 commits. Moreover, for each rule violation, we report the trend of smell oc-15 currences and fixes over time, along with a summary table that describes the 16 most fixed smells. We also discuss interesting cases of smell-fixing commits. 17

## <sup>18</sup> 4.2 RQ<sub>2</sub>: Which Dockerfile smells are developers willing to address?

<sup>19</sup> To answer RQ<sub>2</sub>, we first defined a list of rules, based both on the literature and <sup>20</sup> the results of RQ<sub>1</sub>, and then implemented a rule-based refactoring tool, DOCK-<sup>21</sup> LEANER, to automatically fix them. We defined the fixing rules as described <sup>22</sup> in the *hadolint* documentation<sup>12</sup>. Next, we use DOCKLEANER to fix smells in <sup>23</sup> existing Dockerfiles from open-source projects and submit the changes to the <sup>1</sup> developers through pull requests to understand if they agree with the fixes and <sup>2</sup> are keen to accept them. We describe these steps in the following sections.

<sup>&</sup>lt;sup>12</sup> https://github.com/hadolint/hadolint/wiki

#### 3 4.2.1 Fixing rules for Dockerfile Smells

As a preliminary step, we identified a set of Dockerfile smells that we wanted 4 to fix, considering the list of the most occurring Dockerfile smells, ordered by 5 prevalence, according to the most recent paper on this topic [9]. However, we 6 excluded and added some rule violations. Specifically, among the missing version pinning violations, we excluded DL3013 (*Pin versions in pip*) and DL3018 8 (*Pin versions in apk add*) because they are less occurring variants (*i.e.*, 4%9 and 5%, respectively) of the more prevalent smell DL3008 (15%), even if con-10 cerning different package managers. Additionally, we include in DOCKLEANER 11 the most occurring smells resulting from the analysis performed in  $RQ_1$  and 12 not reported in the literature. We report in Table 1 the full list of smells target 13 in our study, along with the rule we use to automatically produce a fix. It is 14 clear that most of the smells are trivial to fix. For example, to fix the violation 15 DL3020, it is just necessary to replace the instruction ADD with COPY for files 16 and folders. In the case of the version pinning-related smells (i.e., DL3006 and 17 DL3008), instead, a more sophisticated fixing procedure is required. We refer 18 to version pinning-related smells as to the smells related to missing versioning 19 of dependencies and packages. Such smells can have an impact on the repro-20 ducibility of the build since different versions might be used if the build occurs 21 at different times, leading to different execution environments for the applica-22 tion. For example, when the version tag is missing from the FROM instruction 23 of a Dockerfile (*i.e.*, DL3006), the most recent image having the latest tag is 24 automatically selected. To fix such smells, we use a two-step approach: (i) we 25 identify the correct versions to pin for each artifact (e.g., each package), and 26 (ii) we insert the selected versions to the corresponding instruction lines in 27 the Dockerfile. We describe below in more detail the procedure we defined for 28 each smell. 29 **Image version tag (DL3006).** This rule violation identifies a Dockerfile 30 where the base image used in the FROM instruction is not pinned with an explicit 31 tag. In this case, we use a fixing strategy that is inspired by the approach of 32 Kitajima et al. [11]. Specifically, to determine the correct image tag, we use the 33 image name together with the image digest. Docker images are labeled with one 34

or more *tags*, mainly assigned by developers, identifying a specific version of the
 image when *pulled* from DockerHub. On the other hand, the *digest* is a hash
 value that uniquely identifies a Docker image having a specific composition

<sup>38</sup> of dependencies and configurations, automatically created at build time. The <sup>39</sup> *digest* of existing images can be obtained via the DockerHub APIs<sup>13</sup>. Thus,

<sup>40</sup> the only way to uniquely identify an image is using the *digest*. To fix the smell,

<sup>41</sup> we obtain (i) the *digest* of the input Docker image through build, (ii) we find

<sup>42</sup> the corresponding image and its tags using the DockerHub APIs, and (iii) we

<sup>1</sup> pick the most recent tag assigned, that is different from the *"latest"* tag. An

 $_{\rm 2}$   $\,$  example of smell fixed through this rule is reported in Fig. 3.

 $<sup>^{13}</sup>$  https://docs.docker.com/docker-hub/api/latest/

FROM ubuntu	FROM ubuntu:20.04	
(A) Smelly line	(B) Possible solution.	

Fig. 3: Example of rule DL3006.

Pin versions in package manager (DL3008). The version pinning 3 smell also affects package managers for software dependencies and packages 4 (e.g., apt, apk, pip). In that case, differently from the base image, the pack-5 age version must be searched in the source repository of the installed pack-6 ages. The smell regards the *apt* package manager, *i.e.*, it might affect only the 7 Debian-based Docker images. For the fix, we consider only the Ubuntu-based 8 images since (i) we needed to select a specific distribution to handle versions (more on this later), and (ii) Ubuntu is the most widespread derivative of 10 Debian in Docker images [9]. The strategy we use to solve DL3008 works as 11 follows: First, a parser finds the instruction lines where there is the apt com-12 mand, and it collects all the packages that need to be pinned. Next, for each 13 package, the current latest version number is selected considering the OS dis-14 tribution (e.g., Ubuntu, Xubuntu, etc.), and the distro series (e.g., 20.04 Focal 15 Fossa or 14.04 Trusty Tahr). The series of the OS is particularly important, 16 because they may offer different versions for the same package. For instance, 17 if we consider the curl package, we can have the version 7.68.0-1ubuntu2.5 18 for the Focal Fossa series of Ubuntu, while for the series Trusty Tahr it equals 19 to 7.35.0-1ubuntu2.20. So, if we try to use the first in a Dockerfile using 20 the Trusty Tahr series, the build most probably fails. The final step consists 21 in testing the chosen package version. Generally, a package version adopts 22 semantic versioning, characterized by a sequence of numbers in the format 23 <MAJOR>.<PATCH>. However, the specific versions of the packages 24 might disappear in time from the Ubuntu central repository, thus leading to 25 errors while installing them. Given that the PATCH release does not drastically 26 change the functionalities of the package and that old patches frequently dis-27 appear, we replace it with the symbol '\*', indicating "any version," in such a 28 way the *latest* version is automatically selected. After that, a simulation of the 29 apt-get install command with the pinned version is executed to verify that 30 the selected package version is available. If it is, the package can be pinned 31 with that version; otherwise, also the MINOR part of the version is replaced 32 with the '\*' symbol. If the package can still not be retrieved, we do not pin 33 the package, *i.e.*, we do not fix the smell. Pinning a different MAJOR version, 34 indeed, could introduce compatibility issues and the developer should be fully 35 aware of this change. An example of a fix generated through this strategy is 36 reported in Fig. 4. It is worth saying that we apply our fixing heuristic only 37 to packages having missing version pinning. This means that we do not up-38 date packages pinned with another version (e.g., older than the reference date1 used to fix the smell). Moreover, in some cases, developers might not want the 2

RUN apt-get install -y curl	RUN apt-get install -y curl=7.*
(A) Smelly line	(B) Possible solution.

Fig. 4: Example of rule DL3008.

H i! The Dockerfile placed at {dockerfile_path} contains the best practice violation {violation_id} detected by the hadolint tool. The smell {violation_id} occurs when {violation_description}
This pull request proposes a fix for that smell generated by my fixing tool. The patch was manually verified before opening the pull request. To fix this smell, specifically, { <i>fixing_rule_explanation</i> }.
This change is only aimed at fixing that specific smell. If the fix is not valid or useful, please briefly indicate the reason and suggestions for possible improve- ments

Thanks in advance.

Fig. 5: Example of the pull request message. The placeholders (wrapped in curly braces) will be replaced with the corresponding values.

<sup>3</sup> pinned package version, but rather a different one, despite the version we pin

<sup>4</sup> is most likely the closest one to the one they originally tested their Dockerfile

5 on. For example, they want a newer version of that package (e.g., the latest).

<sup>6</sup> We discuss those cases during the evaluation phase of the automated fixes via

7 pull requests.

## 8 4.2.2 Evaluation of Automated Fixes

To evaluate if the fixes generated by DOCKLEANER are helpful, we propose 9 them to developers by submitting the patches on GitHub via pull requests. 10 The first step is to select the most active repositories to ensure responses for 11 our pull requests. To achieve this, we select a subset of repositories from our 12 study context ensuring that, each repository, (i) contains at least one Dockerfile 13 affected by one or more smells that we can fix automatically (reported in 14 Table 1), and (ii) at least one pull request merged, along with commit activity, 15 in the last three months. In this way, we select a total of 186 repositories 16 containing 829 unique Dockerfiles affected by 5,403 smells. The next step is 17 to associate each repository with a specific smell corresponding to a single 18 Dockerfile to fix. This is to avoid flooding developers with pull requests. 19

We used a greedy algorithm to select the smell to fix in the Dockerfiles from the candidate repositories to ensure each of them is considered a bal-

<sup>2</sup> anced number of times. We start from the less occurring smells among all the

<sup>3</sup> available repositories, and we iteratively (i) select one target smell to fix, (ii)

randomly select one Dockerfile candidate containing that smell, (iii) assign the

<sup>5</sup> repository to that smell to mark it as unavailable for the successive iterations,

<sup>6</sup> and (iv) increment a counter, for each smell, of the assigned Dockerfile candi-

<sup>7</sup> dates. The algorithm stops when there are no more repositories available. The

<sup>8</sup> counter of assigned smells is used, along with the overall smell occurrence, in

<sup>9</sup> the first step of the heuristic. This ensures that, for each iteration, we consider <sup>10</sup> the smell (i) having the lower occurrence and (ii) is currently assigned for

the fix to a lower number of repositories. In this phase we manually discard
smells that can not be fixed by DOCKLEANER. For example, for DL3008, we
only support Ubuntu-based Dockerfiles, but the smell might also affect the

<sup>14</sup> Debian-based ones. In total, we excluded 14 smells.

At the end of that procedure, we followed the commonly used git work-15 flow best practices for opening the pull requests. Specifically, we first created 16 a fork for the target repository. Then, we created a branch where the name 17 follows the format fix/dockerfile-smell-DLXXXX. Finally, we signed-off the 18 patches as it is required by some repositories (as well as being a good practice), 19 and we submitted the pull request. To do this, we defined and used a struc-20 tured template for all the pull requests, as reported in Fig. 5. We manually 21 modified the template in the cases where the repository requires a custom-22 defined guidelines. The time required by DOCKLEANER to generate the fixing 23 recommendations is only a few seconds for the simpler fixing procedures (e.g.,24 replacing COPY with ADD). For the more complex ones, such as version pinning, 25 it can even take a few minutes. 26

For the evaluation, we adopted a methodology similar to the one used by 27 Vassallo et al. [20]. In detail, we monitored the status of each pull request for 28 more than 7 months (*i.e.*, 218 days, starting from the last created pull request 29 date) to allow developers to evaluate it and give a response. We interacted 30 with them if they asked questions or requested additional information, but 31 we did not make modifications to the source code of the proposed fix unless 32 they are strictly related to the smell (e.g., the fixing procedure of the smell 33 is reported as not valid). We report such cases in the discussion section. At 34 the end of the monitoring period, we tagged each pull request with one of the 35 following states: 36

<sup>37</sup> – *Ignored*: The pull request does not receive a response;

<sup>38</sup> - *Rejected/Closed*: The pull request has been closed or is explicitly rejected;

<sup>39</sup> – *Pending*: The pull request has been discussed but is still open;

40 - Accepted: The pull request is accepted to be merged but is not merged yet;

 $_{41}$  – *Merged*: The proposed fix is in the main branch.

For each type of fixed smell, we report the number and percentage of the fix recommendations accepted and rejected, along with the rationale in case of rejection and the response time. Also, we conducted a qualitative analysis of the developers' interactions. In particular, we analyzed those where the pull request is rejected or pending to understand why the fix was not accepted. For example, the fix might have been accepted because the developers were



Fig. 6: Occurrence over time for the top 10 Dockerfile smells.

- <sup>3</sup> not interested in performing that modification to their Dockerfile. Moreover,
- <sup>4</sup> we analyze the additional information that the developer submits on rejected
- <sup>5</sup> pull requests, from which we extract takeaways useful for both practitioners
- <sup>6</sup> and researchers. Using a card-sorting-inspired approach [18] performed by two
- $_{7}$  of the authors on the obtained responses, we identified a set of categories that
- $_{\circ}$   $\,$  we used to classify the developers' reactions to rejected pull requests.
- 9 4.2.3 Data Availability
- <sup>10</sup> The code and data used in our study, along with the implementation of DOCK-<sup>11</sup> LEANER, can be found in the replication package [17].

# 12 5 Analysis of the Results

In this section, we report the analysis of the results achieved in our study in
 order to answer our research questions.

# <sup>15</sup> 5.1 RQ<sub>1</sub>: How do developers fix Dockerfile smells?

We report in Fig. 6 the trend of the 10 most occurring Dockerfile smells among the Dockerfile snapshots we analyzed. To plot this figure, we collected all the unique Dockerfiles (based on their path and repository) for each year, then we extracted and counted all the smells of the latest version of each of them (for each year).



Fig. 7: Fixing trend over time for the 10 most fixed Dockerfile smells.



Fig. 8: Overall fixing time delta (days) among all Dockerfile smells.

<sup>3</sup> The most occurring smell is DL3006 – version pinning for the base image–,

<sup>4</sup> followed by DL3008 – missing version pinning for apt-get-, which is also the

 $_{\rm 5}$  most growing one, and DL4000 – deprecated MAINTAINER. Since smell DL4000

 $_{6}$  became a bad practice in 2017<sup>14</sup> after the deprecation of the MAINTAINER

<sup>1</sup> instruction, we excluded its occurrences before that date from the plot.

<sup>&</sup>lt;sup>14</sup> https://docs.docker.com/engine/release-notes/prior-releases/ #1130-2017-01-18



Fig. 9: Cumulative fixes over time interval (days) among all Dockerfile smells.

Table 2: Summary of fixed Dockerfile smells, reporting the number of fixes (manually validated), median time to fix (in days), and the magnitude of changes performed in the repository until the smell has been fixed (median number of commits). Only smells with at least 5 manually validated fixes are reported.

Rule	Description	# Solved	Days (Med.)	Changes (Med.)
DL3059	Consider consolidation for multiple consecutive <b>RUN</b> instructions	168	8.9	4.0
DL3006	Missing version pinning for base image	53	13.7	8.0
DL3007	Avoid to use the latest to tag the version of an image	45	64.6	43.0
DL4000	MAINTAINER is deprecated	45	13.5	1.0
DL3020	Use COPY instead of ADD for files and folders	43	3.8	5.0
DL3003	Use WORKDIR to switch to a directory	29	0.2	0.0
DL3015	Avoid additional packages by specifyingno-install-recommends	26	12.6	4.0
DL3009	Delete the apt-get lists after installing packages	25	23.9	8.0
DL3025	Use arguments JSON notation for CMD and ENTRYPOINT	21	6.0	2.0
DL3048	Invalid Label Key	18	0.1	0.0
DL3019	Use theno-cache switch when installing packages using apk	15	88.0	4.0
DL3004	Do not use sudo as it leads to unpredictable behavior	11	0.3	2.0
DL3028	Pin versions in gem install	8	41.0	57.0
DL3042	Avoid cache directory with pip installno-cache-dir <package></package>	7	0.0	0.0
DL3032	yum clean all missing after yum command	6	597.8	10.0
DL3016	Pin versions in npm	5	33.6	4.0

In our manual validation, we found that 33.6% of the commits in which 2 smells disappear actually fix smells. We report in Table 2 a summary of the 3 characteristics of such commits for the smells for which we found at least 5 4 fixes (from a total of 572 fixed smells). In detail, we report the total number 5 of fixing commits, and the average fixing time, measured both as days and 6 the number of commits that elapsed between the last commit introducing 7 a smell and the smell-fixing commit. Additionally, we report in Fig. 8 the 8 adjusted boxplots describing the days that passed after each smell got fixed. 9 We report in Fig. 7 the fixing trend over time for the 10 most fixed Dockerfile 10 smells. Also, in this case, we consider only the changes which we manually 1 validated as smell-fixing commits. However, this time, we consider each smell 2 fixed separately. This means that, if a commit fixes 5 smells, we count the 3 commit as 5 different fixes, one for each smell. The most fixed smell is DL3059 4

- multiple consecutive RUN instructions. It is worth noting that we found this

fix  $\sim 3$  times more frequently than any other fix. This is because we found

 $_{7}$   $\,$  that, when there are many consecutive RUN instructions, developers tend to fix

<sup>8</sup> all of the occurrences of this issue in a single commit. Other common fixes are

<sup>9</sup> version pinning for base images (DL3006 and DL3007), along with DL4000 <sup>10</sup> – deprecated MAINTAINER and DL3020 – prefer COPY over ADD for files and

10 – deprecated11 folders.

5

We report in Fig. 9 the results of our survivability analysis of the smells 12 by plotting the number of fixed smells in different amounts of time (the time 13 is on a logarithmic scale). It is clear that most of the fixes have been per-14 formed within 1 day (203 instances). This means that when developers intro-15 duce Dockerfile smells, they immediately perform maintenance during the first 16 adoptions. On the other hand, if a smell survives the first day, it is less likely 17 that it gets fixed later. In fact, according to Table 2, the smells that survive 18 the less are DL3048 (incorrect LABEL format) and DL3042 (--no-cache-dir 19 for pip install), which have been fixed in less than one day in most of the 20 cases (100% and 60%, respectively). It is interesting to notice that two similar 21 smells, *i.e.*, DL3006 and DL3007, have largely different survivability. When the 22 latest tag is explicitly used (DL3007) instead of being inferred (DL3006), the 23 smell survives  $\sim 5$  times more (both in terms of days and commits, as reported 24 in Table 2). However, it is worth noting that the effects of both tags are exactly 25 the same. 26 We evaluated how many smell-fixing commits can be considered *informed*.

27 We consider an informed fix when the developer explicitly mentions that the 28 aim of the fix is to remove bad patterns in the commit message. We found that 29 only 18 out of 336 manually validated fixes are *informed*. The most common 30 smell explicitly addressed by developers is DL4000 (fixed in 4 cases) – dep-31 recated MAINTAINER. An example can be found in commit 811582f, from the 32 repository webbertakken/K8sSymfonyReact<sup>15</sup>. Among the remaining ones, 33 DL3025 - JSON notation for CMD and ENTRYPOINT- (4 cases) and DL3020 34 prefer COPY over ADD for files and folders- (3 cases) are the smells of which 35 developers are more aware. 36

As for the *non-informed* cases, mainly developers report that the fix is aimed at (generically) improving the performance of the Dockerfile. Examples are the fixes for rule DL3059 explicitly performed to reduce the Docker image size<sup>16</sup> and the number of layers<sup>17</sup>. In some cases, we found that developers use linters to detect bad practices. Among those, only one commit explicitly mentioned *hadolint*<sup>18</sup>, while in other cases they mentioned the tool *DevOps*-*Bash-tools*<sup>19</sup>.

In the end, we can conclude that developers have a limited knowledge about
 Dockerfile best practices, in terms of the quality of the Dockerfile code. This is

 $^{15}$  https://github.com/webbertakken/K8sSymfonyReact/commit/811582f

<sup>17</sup> https://github.com/Eadom/ctf\_xinetd/commit/21f2785

<sup>18</sup> https://github.com/flyway/flyway-docker/commit/3eeabe5

<sup>19</sup> https://github.com/HariSekhon/Dockerfiles/commit/eeab92a

<sup>&</sup>lt;sup>16</sup> https://github.com/KDE/kaffeine/commit/d03145b

Rule	Ignored	Rejected	Pending	Accepted	Merged*	Assigned
DL4000	1 (8%)	0	0	0	12 (92%)	13
DL3020	2(14%)	2(14%)	0	0	10 (71%)	14
DL3006	2(23%)	2(8%)	0	2(21%)	9 (69%)	13
DL3007	2(14%)	3(21%)	0	0	9 (64%)	14
DL3015	3(23%)	2(15%)	0	0	8 (62%)	13
DL3025	4 (33%)	0	0	1(8%)	8(67%)	12
DL3059	2(15%)	3(23%)	0	0	8(62%)	13
DL3048	2(22%)	1(11%)	0	0	6(67%)	9
DL3009	5(42%)	3(25%)	0	2(17%)	4 (33%)	12
DL3003	1 (14%)	3(43%)	0	0	3(43%)	7
DL3008	5(33%)	7(47%)	0	0	3(20%)	15
DL4006	4(50%)	1(13%)	0	0	3(38%)	8
Total	33 (23%)	27 (19%)	0	5(3%)	83 (58%)	143

Table 3: Opened pull requests and their resulting status sorted by number of accepted and merged PRs. The column *Merged*<sup>\*</sup> reports the cumulative number of accepted patches (sum of accepted and merged).

<sup>4</sup> because they are more interested in the optimization of other non-functional

 $_{\rm 5}$   $\,$  aspects such as build time and size of the Docker image.

**Q** Summary of  $\mathbf{RQ}_1$ : The most fixed smells are those related to consecutive RUN instructions (DL3059), version pinning for the base image (DL3006/DL3007), use of the deprecated MAINTAINER instruction (DL4000) along with the usage of WORKDIR to change directory (DL3020). The 34% of the evaluated commits (1000) actually fixed the smell. Also, most of the smells are fixed immediately after their introduction (within 1 day) and, when this does not happen, they might remain in the repository for a long time (more than 3 years).

# 7 5.2 RQ<sub>2</sub>: Which Dockerfile smells are developers willing to address?

In Table 3 we report the results of the evaluation performed via GitHub pull 8 requests. In total, we submitted 143 pull requests. The majority of them have 9 been accepted or merged by developers (58%). On the other hand, 23% them 10 have been ignored, while 19% received an explicit rejection from the developers. 11 The smells receiving the highest acceptance rate are DL4000 – deprecated 12 MAINTAINER- (92%) and DL3020 - prefer COPY over ADD for files and folders-1 (71%), followed by rule DL3006 – version pinning for the base image- (69%). 2 This is similar to what we reported for RQ<sub>1</sub>, where they resulted to be the most 3 fixed smell among the manually validated smell-fixing commits. This means 4 that developers care about those smells as they frequently fixed them and they 5 are also willing to accept fixes. The smell DL3008 – missing version pinning 6 for apt-get- has been the most rejected fix (47% acceptance), with only 3

6



Fig. 10: Average resolution time (days) for merged pull requests (a) and rejected pull requests (b).

accepted pull requests, along with smell DL4006 - use of pipefail for piped
operations- which has been the most ignored one (50%). The low acceptance
rate (33%) resulting for smell DL3009 (deletion of apt-get sources lists) is
surprising, since developers are prone to reduce the image size, as we noticed in
RQ<sub>1</sub>. Despite this, we can conclude that they do not prefer to remove apt-get
source lists to achieve this goal.
In Fig. 11 we report the adjusted boxplot for the time required for pull

requests to get the first response and to be resolved. Additionally, Fig. 10
reports the median resolution time, measured in days, of the submitted pull
requests by smell type. For both of those figures, we only consider merged
and rejected PRs, because they are the ones for which we have a definitive



Fig. 11: Adjusted boxplot of the number of days required for a pull request to obtain a response (left) and to be merged/rejected (right).

# response from the developers. The smell DL3025 – JSON notation for CMD and ENTRYPOINT- is the one that has been accepted in the shortest time interval, followed by DL3006 – version pinning for the base image- and DL4000 – deprecated MAINTAINER. Despite the fixes for DL3020 – prefer COPY over ADD for files and folders- are the second most-accepted ones, they have a median of 5 days to get accepted and merged.

On the other hand, the fixes for DL4006 – use of pipefail for piped operations- have been rejected almost immediately by developers. This also happens for DL3008 – missing version pinning for apt-get.

Finally, we report in Table 4 the reasons why developers rejected our pull requests. We assigned one or more categories, for each rejected change, by analyzing the responses for the 27 rejected pull requests. Most of the time, the fix has been considered invalid (22% of cases). This means that the proposed change was not a valid improvement for the Dockerfile. In 11% of cases, the developers did not accept the change as they use the Dockerfile in testing or development environments.

The rejections of the fixes for DL3008 are interesting: In 19% of the cases, the changes have been rejected because they are not perceived as a concrete fix. Furthermore, the fixes for that smell have been rejected because they could negatively impact the security of the image (8% of cases) or cause a build failure in the future (4% of cases).

Reason	Involved Smells	Occurrences
Invalid fix	DL3003,DL3007,DL3020,DL3059,DL4006	6
No reason	DL3006,DL3008,DL3015,DL3059	4
Fix not required	DL3008,DL3059	5
Not trusted	DL3007,DL3008,DL3020	3
Testing environment	DL3006,DL3007,DL3008	3
Reduces security	DL3008	2
Development environment	DL3009	2
Vendored dependency	DL3003	1
Potential breaking change	DL3008	1
Unused file	DL3009	1

Table 4: Categories of reasons why developers rejected our pull requests.

**Q** Summary of RQ<sub>2</sub>: Developers accepted most of the Dockerfile smell fixes we provided (58%) and rejected only a few of them (19%). They particularly liked the fixes for DL4000 (deprecated MAINTAINER), DL3020 (prefer COPY over ADD for files and folders), and DL3006 (version pinning for the base image). Instead, they frequently rejected DL3008 (version pinning for apt-get packages) (47%). The reason is that it is seen as a bad practice as it could lead to failures or security issues in the future.

#### 4 6 Discussion

Despite the majority of the submitted pull requests got accepted, there are 5 some specific smells that developers are not willing to address. Looking at Table 4, in 5 cases, the fix was rejected because the container was used in a testing 7 or development environment. An example is the fix proposed for  $DL3009^{20}$ , 8 where, even if the change can reduce the image size, it negatively impacts the 9 image build time. Thus, for that reason, the change has been rejected. Prob-10 ably, the concern about build time comes from frequent builds performed for 11 that specific Dockerfile. A different example is the pull request submitted to 12 envoyproxy/ratelimit<sup>21</sup>, the reason for the rejection is that developers do 13 not care about the version pinning (DL3007) as they use that Dockerfile for 14 testing and they need to test the latest version of the software. This is not the 15 same for DL3006 when the tag is missing. In that case, developers are more 16 likely to accept the version pinning for the base image (see RQ<sub>1</sub> and RQ<sub>2</sub>). 17

<sup>&</sup>lt;sup>20</sup> https://github.com/Shopify/semian/pull/484

 $<sup>^{21}\ {\</sup>tt https://github.com/envoyproxy/ratelimit/pull/411}$ 

**Q** Lesson 1. Developers tend to use the "latest" tag for the base images (DL3007) in order to obtain the latest version of the image, while they are willing to accept the version pinning when the tag is missing (DL3006). However, as the "latest" tag is not immutable, this practice can lead to unexpected behaviors when the base image is updated.

DL3008 constitutes a peculiar case. Fixing such a smell requires devel-1 opers to pin the version of the apt-get packages to make the build more 2 reproducible. Developers, however, believe that doing so might be mislead-3  $ing^{22}$ , or it might make the build more fragile<sup>23</sup>. Indeed, this happened for 4 an accepted pull request, where after a month the version pinning for the 5 package ca-certificates caused a build failure because the pinned version 6 was not available anymore<sup>24</sup>. Moreover, the smell DL3008 led to interesting 7 discussions. For example, a suggestion was to provide an automated script 8 to periodically pin the package versions when there is an  $update^{25}$ . For 3 of 9 the proposed fixes, the developers additionally highlighted that they do not 10 trust the change because it has been generated by an automated tool. This 11 happened even if we specified that we manually checked the correctness of the 12 change. 13

**Q** Lesson 2. Version pinning for OS packages is not considered a good practice. Developers tend to avoid it because (i) they consider it a misleading practice, (ii) it could lead to building failures due to the unavailability of the pinned version, and (iii) missed security updates when the pinned version gets older.

In 6 cases, instead, developers did not perceive the change as correct or 15 sufficient for a fix. This happens, for example, in commits 5531f2e<sup>26</sup> (DL3020) 16 and 320ba87<sup>27</sup> (DL4006). An interesting discussion arose for the rejected fix 17 of DL3003<sup>28</sup>. The fix for that smell provides the replacement of "cd <path>" 18 with "WORKDIR <path>". However, for that particular case, fixing the smell 19 required putting a WORKDIR instruction before the smelly code block and an-20 other after to switch back to the previous working directory. This is because 21 the target smelly code temporarily changes the working directory to operate 22 on specific files. In other words, there are cases in which developers believe it 23 is legitimate to change the working directory through cd (mostly, when this 24 change is temporary). We report an example in Fig. 12, where the fix has been 25 rejected because the change of the working directory is temporary. 26

22

18

14

 $<sup>^{22}</sup>$  https://github.com/James-Yu/LaTeX-Workshop/pull/3837

<sup>&</sup>lt;sup>23</sup> https://github.com/Yelp/aactivator/pull/47

 $<sup>^{24} \ {\</sup>tt https://github.com/FDio/govpp/pull/123}$ 

 $<sup>^{25}\ {\</sup>tt https://github.com/Lookyloo/lookyloo/pull/663}$ 

 $<sup>^{26} \ \</sup>texttt{https://github.com/ROCmSoftwarePlatform/Tensile/pull/1707}$ 

<sup>&</sup>lt;sup>27</sup> https://github.com/bupy7/xml-constructor/pull/6

<sup>28</sup> https://github.com/NUbots/NUbots/pull/1063

		@@ -53,7 +53,8 @@ COPY etc/ld.so.conf.d/usrlocal.conf /etc/ld.so.conf.d/usrlocal.conf
53	53	RUN ldconfig
54	54	
55	55	# Make a symlink from /usr/local/lib to /usr/local/lib64 so library install location is irrelevant
56		- RUN cd /usr/local && ln -sf lib lib64
	56	+ WORKDIR /usr/local
	57	+ RUN ln -sf lib lib64
57	58	
58	59	# Generate toolchain files for the generic platform
59	60	<pre>COPY usr/local/toolchain/generate_toolchains.py /usr/local/generate_toolchains.py</pre>

Fig. 12: Example of a wrong fix for DL3003. In that case, the change of working directory is temporary, and the fix has been rejected.

We conclude that, in similar cases, the detected smell is a false positive. 27 This is because the fix will increase the number of layers, in addition to redun-1 dant instructions. This negatively impacts the code quality of the Dockerfile. 2 Comparing the results from  $RQ_1$  and  $RQ_2$ , we can conclude that there 3 are no big differences between the fixes that developers have applied and the 4 changes that we propose via pull requests. The most performed fixes, which 5 are also in the most accepted pull requests, are those related to deprecated 6 MAINTAINER (DL4000), version pinning for the base image (DL3006), and mul-7 tiple consecutive **RUN** instructions (DL3059). There is a difference in terms of 8 the most fixed one. While in the wild developers tend to fix more DL3059, in 9 our pull request the most fixed one is DL4000. As also shown in  $RQ_1$ , they 10 pay more attention to performance improvements over code quality, for which 11 they are not fully aware of what is the current writing best  $practices^{29}$ . In 12 fact, DL4000 is purely related to writing best practices and does not affect 13 performance. When faced with a ready-to-use fix, however, they tend to prefer 14 the ones that more likely will not disrupt the Dockerfile. 15 In general, developers keep more attention to the impact of the change on 16

<sup>17</sup> the build process and the image size, instead of the impact on the quality of <sup>18</sup> the Dockerfile code. Reporting an example among the accepted pull requests, <sup>19</sup> we have the fix proposed for the smell DL3015 (--no-install-recommend <sup>20</sup> flag for apt)<sup>30</sup>, where the developers explicitly asked to fix another Dockerfile <sup>21</sup> affected by the same smell because it decreases the size of the built image.

**Q** Lesson 3. Developers are not fully aware of the best practices for writing Dockerfiles, and they tend to prefer performance improvements over code quality.

Additionally, it is interesting to analyze more in depth the differences in terms of performed fixes for DL3048 (incorrect LABEL format) and DL4000 (MAINTAINER is deprecated, replace with LABEL). Actually, there are two possible ways to format Dockerfile labels. The first one follows the standard format

22

<sup>&</sup>lt;sup>29</sup> https://github.com/riga/law/pull/152

<sup>&</sup>lt;sup>30</sup> https://github.com/lablup/backend.ai/pull/1216

defined by open containers <sup>31</sup>, which is also suggested for DL4000 in the official 27 Docker documentation <sup>32</sup>. The second is more general and does not enforce a pre-defined format. It is reported in the *hadolint* documentation  $^{33}$ , which also 2 is reported in the official Docker documentation as examples of LABEL instruc-3 tions  $^{34}$ . The fixes that we analyzed in RQ<sub>1</sub> that follow the first format are 4 limited only to one repository<sup>35</sup>. In other cases, developers adopted the sec-5 ond format<sup>36</sup>. The fixes proposed via pull requests, instead, follow the second 6 format where for DL4000 we got the highest acceptance rate. This is probably 7 because the second format is more general, avoiding unnecessary constraints 8 and changes on the LABEL instructions $^{37}$ . 9 Moreover, while in this context the fix is still sufficient to correct the smell. 10 in other contexts our fixing procedure could not be correct. The most evident 11 case is for the smell DL3059 (multiple consecutive RUN instructions). In fact, 12 open-source developers tend to fix it mainly by compacting the installation of 13 software packages<sup>38</sup>. In our pull requests, instead, we merge all the subsequent 14 RUN instructions until a comment or a different instruction is found. This could 15 mean that a more complex and informed fixing procedure should be adopted 16 in order to better improve the size and performance of Dockerfiles. Thus, a 17 more advanced approach in that direction could be useful to improve the fixing 18 procedure, taking also into account the aspects that developers are interested 19 to improve (image size and build time). To this aim, considering the scenario 20 in which we are using a **debian** base image, an advanced approach to fix smell 21 DL3059 could be a heuristic that (i) selects all the RUN instructions that are 22 aimed at installing dependencies, (ii) extracts the list of such dependencies, 23 taking also into account if they require external sources lists, and (iii) combine 24 all those installations into a single **RUN** instruction at the top of the Dockerfile. 25 In this way, the re-build time will be reduced thanks to the layers caching 26 system. At the same time, the image size will be reduced since there will be 27 fewer layers and less space wastage (e.g., package cache). For smell DL3003, 28 instead, an advanced fixing approach should target the bash code to correct the 29 usage of the pattern "RUN cd  $\ldots$ ", rather than using WORKDIR. In the example 30 reported in Fig. 12, the smell could be fixed by using the absolute paths instead 31 of the relative paths for the command (e.g., "RUN ln -sf /usr/local/lib 32 /usr/local/lib64"). While this can be done in this case, there are other 33 scenarios in which this could be detrimental. For example, if a custom script 34 writes the output files in the current directory, it is still necessary to use cd 35 before running it. Thus, such a fixing procedure should be applied only for 36 specific bash instructions patterns (like the previously-mentioned one). 37

<sup>&</sup>lt;sup>31</sup> https://specs.opencontainers.org/image-spec/annotations/

 $<sup>^{32}\ {\</sup>tt https://docs.docker.com/engine/reference/builder/{\tt #maintainer-deprecated}}$ 

<sup>&</sup>lt;sup>33</sup> https://github.com/hadolint/hadolint/wiki/DL3048

 $<sup>^{34}</sup>$  https://docs.docker.com/engine/reference/builder/#label

<sup>&</sup>lt;sup>35</sup> https://github.com/HariSekhon/Dockerfiles/commit/f329b94

<sup>&</sup>lt;sup>36</sup> https://github.com/scossu/lakesuperior/commit/a552ff7

<sup>&</sup>lt;sup>37</sup> https://github.com/hpc/charliecloud/pull/1628

<sup>&</sup>lt;sup>38</sup> https://github.com/hpc/charliecloud/commit/aae89d7

**Q** Lesson 4. A more advanced fixing procedure is required for some types of smells (*e.g.*, DL3003 – Use WORKDIR to switch to a directory– and DL3059 – multiple consecutive RUN instructions), *i.e.*, taking into account the context in which the smell is found.

#### <sup>1</sup> 7 Threats to Validity

20

Construct Validity. The threats to construct validity are about the non-2 measurable variables of our study. More specifically, our study is heavily based 3 on the rule violations detected by hadolint. Other tools are able to detect 4 bad practices in Dockerfiles, such as  $dockle^{39}$ . We choose hadolint which is 5 commonly used in the literature [6, 9, 14, 21] and also in enterprise tools for 6 code quality<sup>40</sup>. However, *hadolint* could lead to false positives or can miss some 7 smells<sup>41</sup>. The manual evaluation we performed on the smell-fixing commits 8 validated the identified smells and those that have been removed. During that 9 evaluation, we noticed that hadolint mainly fails to detect the rule DL3059 10 (consecutive RUN instructions). To reduce this impact of this threat on our 11 study, we manually annotated the lines in which the smell was present. 12

Internal Validity. The threats to internal validity are about the design 13 choices that we made which could affect the results of the study. In detail, we 14 used as a study context a sample of repositories extracted from the dataset 15 provided by Eng et al. [9] by considering only those having stargazers count 16 greater or equal to 10. This is commonly used in the literature to avoid toy 17 projects [8]. There can be a bias in the selected smells for our fix recommen-18 dations. We selected the most occurring smell as described in the analysis of 19 Eng et al. [9]. We assume that an automated approach would have the biggest 20 impact on the smells that occur more frequently. Also, at least for some of 21 them, the reason behind the fact that they do not get fixed might be that 22 they are not trivial (*i.e.*, an automated tool would be helpful). The fixing pro-23 cedure for some of the selected smells can be wrong, and some smells might 24 not get fixed. We based the rules on the fixing procedure on the Docker best 25 practices and on the hadolint documentation. Still, to minimize the risk of this, 26 we double-checked the modifications before submitting the pull requests and 27 manually excluded the ones that make the build of the Dockerfile fail. Thus, we 28 ensured the correctness of the fixes generated by DOCKLEANER, submitted via 29 the pull requests, for the cases evaluated in our study. However, it is still pos-30 sible that the tool produces wrong fixes for other Dockerfiles. For example, the 31 version pinning fixes could fail in the cases in which the package is not reach-32 able (*i.e.*, DL3008), or the Docker image digest is not available in DockerHub 33 (*i.e.*, for smells DL3006 and DL3007). It is worth noting, indeed, that our aim 34 is not to evaluate the tool, but rather to understand if developers are willing 35

<sup>&</sup>lt;sup>39</sup> https://github.com/goodwithtech/dockle

<sup>&</sup>lt;sup>40</sup> https://github.com/codacy/codacy-hadolint

<sup>&</sup>lt;sup>41</sup> https://github.com/hadolint/hadolint/issues/693

to accept fixes. Moreover, there is a possible subjectiveness introduced of the 36 manual validation of the smell-fixing commits, which has been mitigated with the involvement of two of the authors and the discussion of the conflicts. Also, 2 it is important to say that the two evaluators have more than 3 years of experi-3 ence with Dockerfiles development and Docker technology in general, allowing 4 them to have a good understanding of the smells and the applied fixes. Finally, 5 we performed the selection of the *last-smell-introducing* commits by using the 6 git blame command on the smelly lines identified by hadolint. Since hadolint 7 can fail to detect some smells, in some cases, the lines impacted by the fix are 8 different from the ones identified by *hadolint*. This means that we got some 9 false positives while we identify the *last-smell-introducing* commits. Since our 10 results showed that Dockerfiles are not frequently changed, we believe that the 11 impact of this threat is limited. 12

External Validity. External validity threats concern the generalizability of our results. In our study, we considered a sample of repositories from GitHub containing only open-source Dockerfiles. This means that our findings might not be generalized to other contexts (*e.g.*, industrial projects) as developers could handle smell in a different way.

## 18 8 Conclusion

In the last few years, containerization technologies have had a significant im-19 pact on the deployment workflow. Best practice violations, namely Dockerfile 20 smells, are widely spread in Dockerfiles [6, 9, 14, 21]. In our empirical study, 21 we evaluated the Dockerfile smell survivability by analyzing the most fixed 22 smells in open-source projects. We found that Dockerfile smells are widely 23 diffused, but developers are becoming more aware of them. Specifically, for 24 those that result in a performance improvement. In addition, we evaluated to 25 what extent developers are willing to accept fixes for the most common smells, 26 automatically generated by a rule-based tool. We found that developers are 27 willing to accept the fixes for the most commonly occurring smells, but they 28 are less likely to accept the fixes for smells related to the version pinning of 29 OS packages. To the best of our knowledge, this is the first in-depth analy-30 sis focused on the fixing of Dockerfile smells. We also provide several lessons 31 learned that could guide future research in this field and help practitioners in 32

- <sup>33</sup> handling Dockerfile smells.
- Acknowledgements The work by Rocco Oliveto, Giovanni Rosa, and Simone Scalabrino was supported by the European Union - NextGenerationEU through the Italian Ministry of University and Research, Projects PRIN 2022 "QualAI: Continuous Quality Improvement
- 37 of AI-based Systems", grant n. 2022B3BP5S, CUP: H53D23003510006.
- The authors would like to thank Alessandro Giagnorio (University of Molise, Italy) for implementing a preliminary version of DOCKLEANER.

#### 40 Conflict of interest

<sup>1</sup> The authors declare that they have no conflict of interest.

#### 2 References

- 3 1. Best practices for writing Dockerfiles. [Online; accessed 2-Jun-2022]
- hadolint: Dockerfile linter, validate inline bash, written in Haskell. [Online; accessed
   28-May-2022]
- 6 3. ShellCheck, a static analysis tool for shell scripts. [Online; accessed 2-Jun-2022]
- Azuma, H., Matsumoto, S., Kamei, Y., Kusumoto, S.: An empirical study on selfadmitted technical debt in dockerfiles. Empirical Software Engineering 27(2), 1–26 (2022)
- 5. Becker, P., Fowler, M., Beck, K., Brant, J., Opdyke, W., Roberts, D.: Refactoring:
   Improving the Design of Existing Code. Addison-Wesley Professional (1999)
- Cito, J., Schermann, G., Wittern, J.E., Leitner, P., Zumberi, S., Gall, H.C.: An empirical analysis of the docker container ecosystem on github. In: 2017 IEEE/ACM 14th
   International Conference on Mining Software Repositories (MSR), pp. 323–333. IEEE (2017)
- 7. Cohen, J.: A coefficient of agreement for nominal scales. Educational and psychological
   measurement 20(1), 37-46 (1960)
- Babic, O., Aghajani, E., Bavota, G.: Sampling projects in github for msr studies.
   In: 2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR), pp. 560–564. IEEE (2021)
- 9. Eng, K., Hindle, A.: Revisiting dockerfiles in open source software over time. In: 2021
   IEEE/ACM 18th International Conference on Mining Software Repositories (MSR), pp.
   449–459. IEEE (2021)
- Henkel, J., Bird, C., Lahiri, S.K., Reps, T.: Learning from, understanding, and supporting devops artifacts for docker. In: 2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE), pp. 38–49. IEEE (2020)
- 11. Kitajima, S., Sekiguchi, A.: Latest image recommendation method for automatic base
   image update in dockerfile. In: International Conference on Service-Oriented Comput ing, pp. 547–562. Springer (2020)
- 12. Kruchten, P., Nord, R.L., Ozkaya, I.: Technical debt: From metaphor to theory and
   practice. Ieee software 29(6), 18–21 (2012)
- 13. Ksontini, E., Kessentini, M., Ferreira, T.d.N., Hassan, F.: Refactorings and technical
   debt in docker projects: An empirical study. In: 2021 36th IEEE/ACM International
   Conference on Automated Software Engineering (ASE), pp. 781–791. IEEE (2021)
- Lin, C., Nadi, S., Khazaei, H.: A large-scale data set and an empirical study of docker
   images hosted on docker hub. In: 2020 IEEE International Conference on Software
   Maintenance and Evolution (ICSME), pp. 371–381. IEEE (2020)
- Ma, Y., Bogart, C., Amreen, S., Zaretzki, R., Mockus, A.: World of code: an infrastructure for mining the universe of open source vcs data. In: 2019 IEEE/ACM 16th
   International Conference on Mining Software Repositories (MSR), pp. 143–154. IEEE
   (2019)
- Regier, D.A., Narrow, W.E., Clarke, D.E., Kraemer, H.C., Kuramoto, S.J., Kuhl, E.A.,
   Kupfer, D.J.: Dsm-5 field trials in the united states and canada, part ii: Test-retest
   reliability of selected categorical diagnoses. American Journal of Psychiatry 170(1),
   59-70 (2013). DOI 10.1176/appi.ajp.2012.12070999. URL https://doi.org/10.1176/
   appi.ajp.2012.12070999. PMID: 23111466
- Rosa, G., Zappone, F., Scalabrino, S., Oliveto, R.: Replication package (2024). https:
   //doi.org/10.6084/m9.figshare.23522679
- 49 18. Spencer, D.: Card sorting: Designing usable categories. Rosenfeld Media (2009)
- Tufano, M., Palomba, F., Bavota, G., Oliveto, R., Di Penta, M., De Lucia, A., Poshyvanyk, D.: When and why your code starts to smell bad (and whether the smells go
- <sup>52</sup> away). IEEE Transactions on Software Engineering **43**(11), 1063–1088 (2017)

- Vassallo, C., Proksch, S., Jancso, A., Gall, H.C., Di Penta, M.: Configuration smells in
  continuous delivery pipelines: a linter and a six-month study on gitlab. In: Proceedings
  of the 28th ACM Joint Meeting on European Software Engineering Conference and
  Symposium on the Foundations of Software Engineering, pp. 327–337 (2020)
- 4 21. Wu, Y., Zhang, Y., Wang, T., Wang, H.: Characterizing the occurrence of dockerfile
- smells in open-source software: An empirical study. IEEE Access 8, 34127–34139 (2020)
  22. Zerouali, A., Mens, T., Robles, G., Gonzalez-Barahona, J.M.: On the relation between outdated docker containers, severity vulnerabilities, and bugs. In: 2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER),
  - pp. 491–501. IEEE (2019)
- Zhang, Y., Vasilescu, B., Wang, H., Filkov, V.: One size does not fit all: an empirical
   study of containerized continuous deployment workflows. In: Proceedings of the 2018
   26th ACM Joint Meeting on European Software Engineering Conference and Sympo-
- sium on the Foundations of Software Engineering, pp. 295–306 (2018)

9



14

2

3

Giovanni Rosa is a Ph.D. student in Software Engineering at University of Molise (UNIMOL), Italy. He received his MSc. in Software System Security, from the same University, in October 2020. His research activity is focused (but not limited to) on software quality and maintenance of DevOps artifacts. More information: https://giovannirosa.com/



Federico Zappone has received his MS degree in Software System Security from the University of Molise (UNI-MOL). He is a co-author of various papers about Distributed Ledger Technologies like Hyperledger Fabric. He is the co-founder of two different companies, Just Another SRL, an Italian company dedicated to developing innovative, high-quality systems, and BB-SMILE SRL, a spin-off born out of a collaboration between the University of Molise and Sapienza University of Rome.



Simone Scalabrino is an Assistant Professor at the University of Molise, Italy, where he leads the DEVelopercentrIc Software Engineering Research group (DEVISER). He has received his PhD degree from the University of Molise in 2019, defending a thesis on automatically assessing and improving source code readability and understandability. His main research interests include code quality, software testing, and empirical software engineering. He is the author of more than 50 papers published in international journals and conferences, and he has received three

ACM SIGSOFT Distinguished Paper Awards at ICPC 2016, ASE 2017, and MSR 2019. He is also co-founder of Datasound srl, a spin-off of the University of Molise. More information: https://dibt.unimol.it/sscalabrino



**Rocco Oliveto** is a Full Professor at the University of Molise (Italy). He is the founder of the Software Engineering and Knowledge Engineering (STAKE) Lab of the University of Molise. Prof. Oliveto is co-author of about 200 papers on topics related to software traceability, software maintenance and evolution, and empirical software engineering. He has received several awards for his research activity, including 5 ACM SIGSOFT Distinguished Paper Awards and 3 Most Influential Paper Awards. Prof. Oliveto participated in the organization and was a mem-

ber of the program committee of several international conferences in the field of software engineering. Since 2018 he has been CEO of Datasound srl, a spinoff of the University of Molise that was created to conceive, design and develop innovative recommendation systems to be applied in different contexts.