# What Quality Aspects Influence the Adoption of Docker Images?
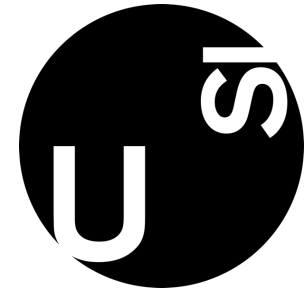
**Giovanni Rosa**, Simone Scalabrino, Gabriele Bavota and Rocco Oliveto

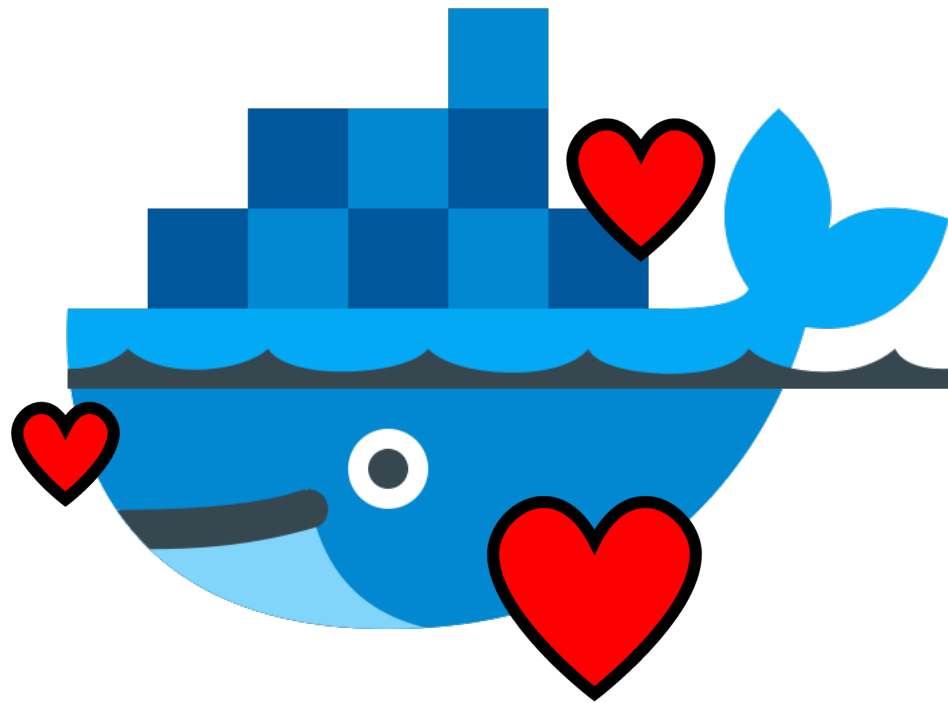@giovannipink

**University of Molise, Italy**

**#1 most-desired**
and
**#1 most-used**
dev tool

2023
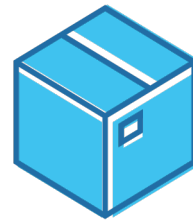Developer
Survey

Why Docker?

**base image**

```
 1 FROM node:12-alpine
 2
 3 RUN apk add --no-cache python2 g++ make
 4
 5 WORKDIR /app
 6 COPY . .
 7
 8 RUN yarn install --production
 9
10 CMD ["node", "src/index.js"]
11
12 EXPOSE 3000 here
```
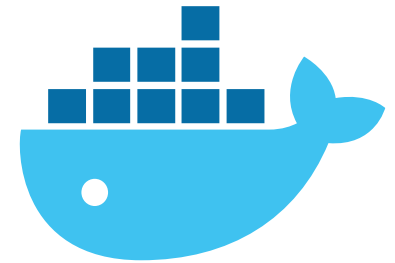
# Dockerfile

base image

```
1  FROM node:12-alpine
2
3  RUN apk add --no-cache python2 g++ make
4
5  WORKDIR /app
6  COPY . .
7
8  RUN yarn install --production
9
10 CMD ["node", "src/index.js"]
11
12 EXPOSE 3000 here
```

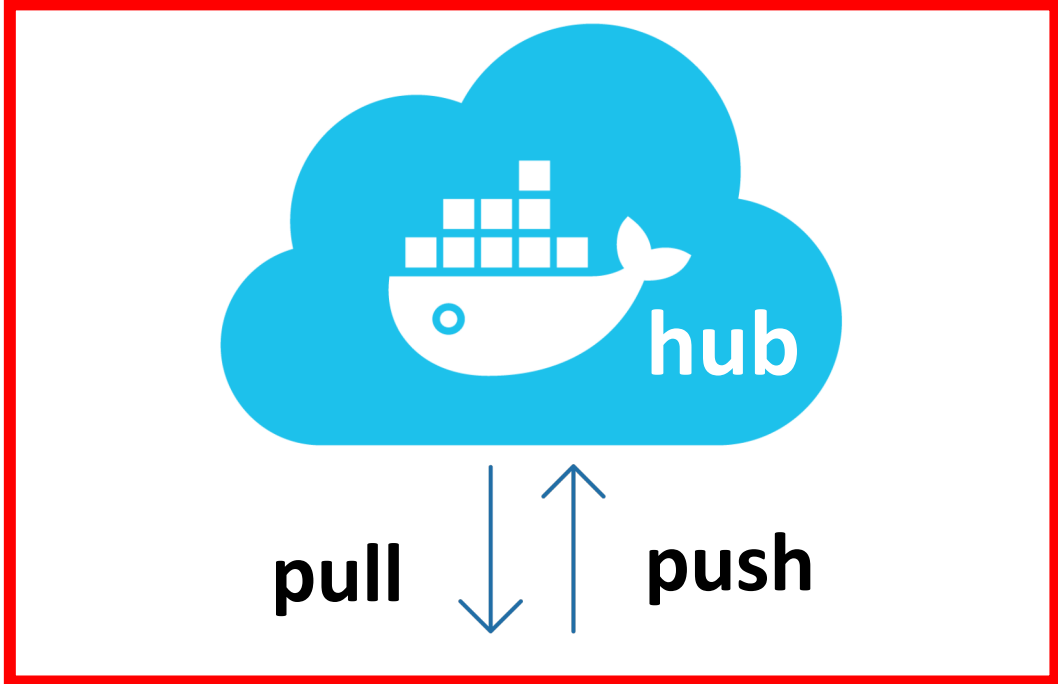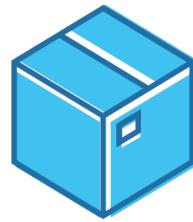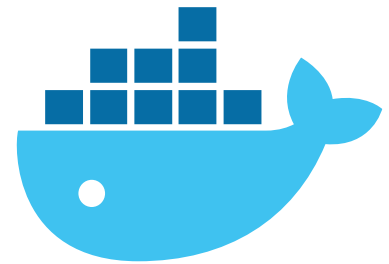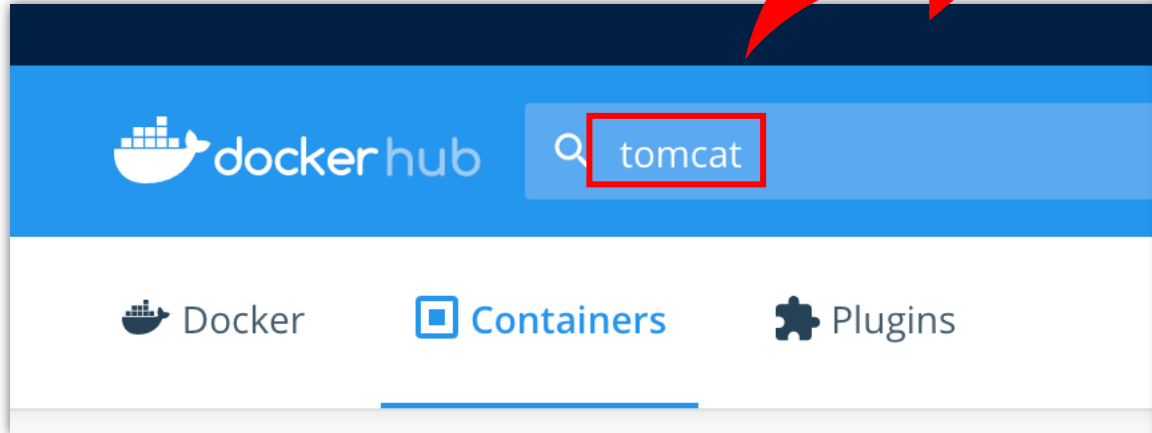Dockerfile → **build** → Image → **run** → Container

Docker in a nutshell

base image

```
1  FROM node:12-alpine
2
3  RUN apk add --no-cache python2 g++ make
4
5  WORKDIR /app
6  COPY . .
7
8  RUN yarn install --production
9
10 CMD ["node", "src/index.js"]
11
12 EXPOSE 3000 here
```

**pull**          **push**

**build**          **run**

Dockerfile          Image          Container

Docker in a nutshell

Which Docker image to choose?

1 - 25 of 10,000 results for **tomcat**.

**tomcat** ✪ **Docker Official Image** · ⬇ 500M+ ·

Updated 3 days ago

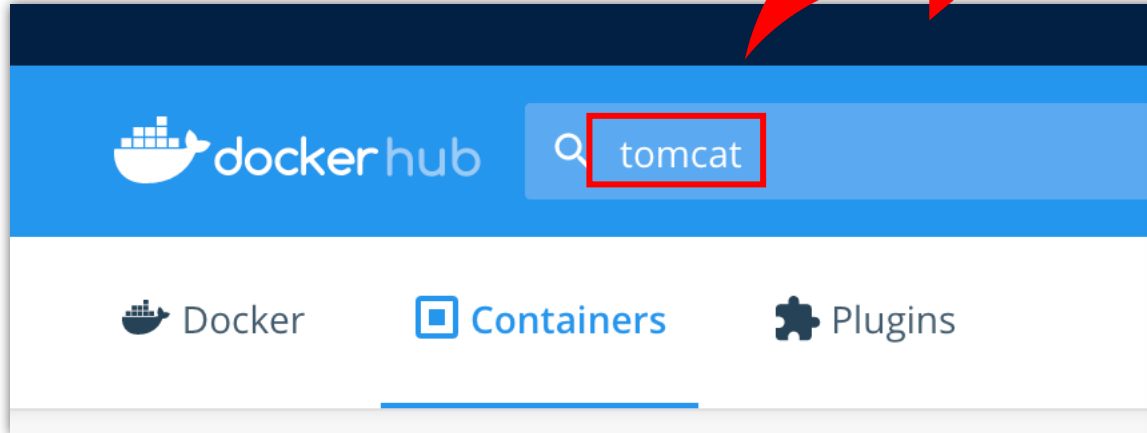Apache Tomcat is an open source implementation
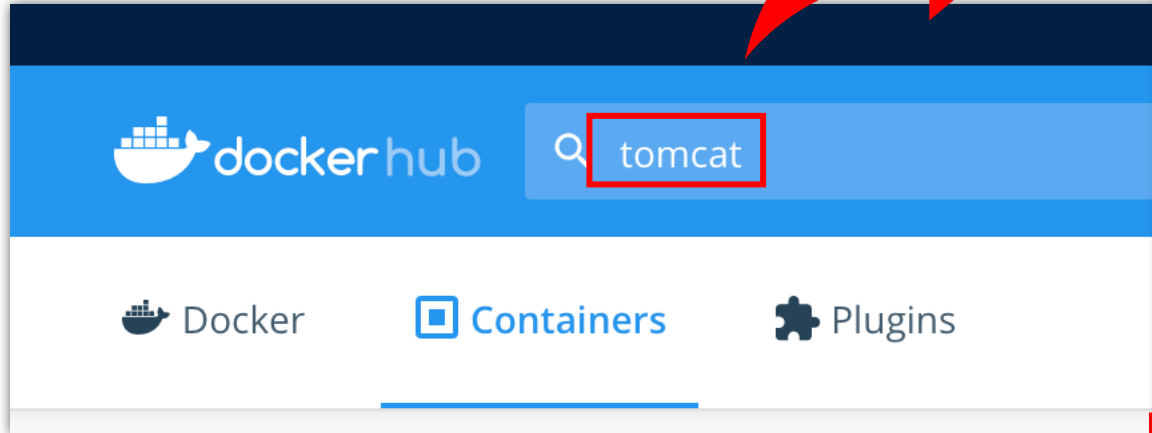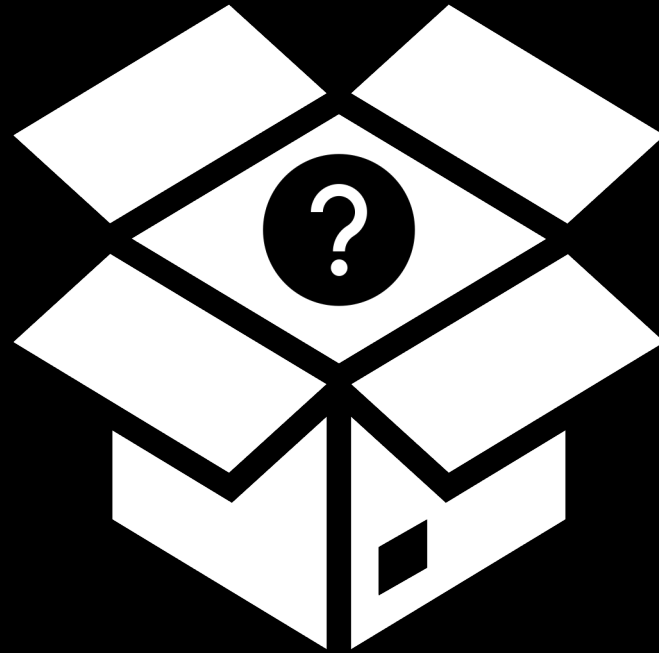
Linux  x86-64  ARM  ARM 64  PowerPC 64 LE  IBM

**jelastic/tomcat** · ⬇ 10M+ · ☆ 4

By jelastic · Updated 6 days ago

An image of the Tomcat Java application server

Which Docker image to choose?

Which Docker image to choose?

How to describe

a «good» Docker image?

"smells are very common in Dockerfile codes"

## Characterizing the Occurrence of Dockerfile Smells in Open-Source Software: An Empirical Study

**Wu et. al 2020**

**Shu et. al 2017**

A Study of Security Vulnerabilities on Docker Hub

"smells are very common in Dockerfile codes"

"images contain more than 180 vulnerabilities on average"

**Wu et. al 2020**

Characterizing the Occurrence of Dockerfile Smells in Open-Source Software: An Empirical Study

**Shu et. al 2017**

A Study of Security Vulnerabilities on Docker Hub

**Ibrahim et. al 2020**

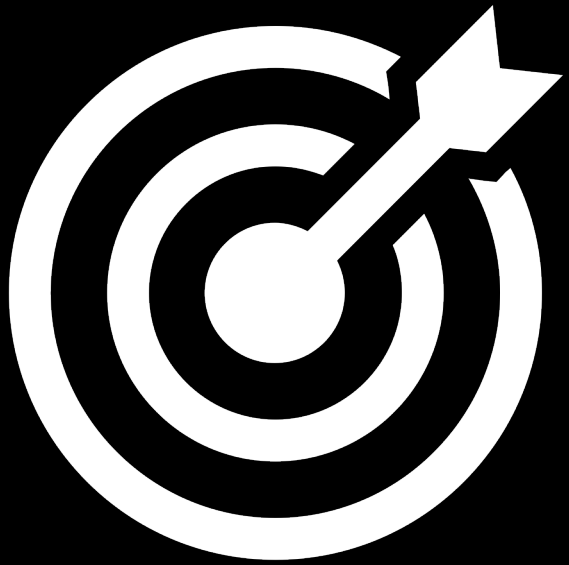Too many images on DockerHub! How different are images for the same system?

"smells are very common in Dockerfile codes"

"images contain more than 180 vulnerabilities on average"

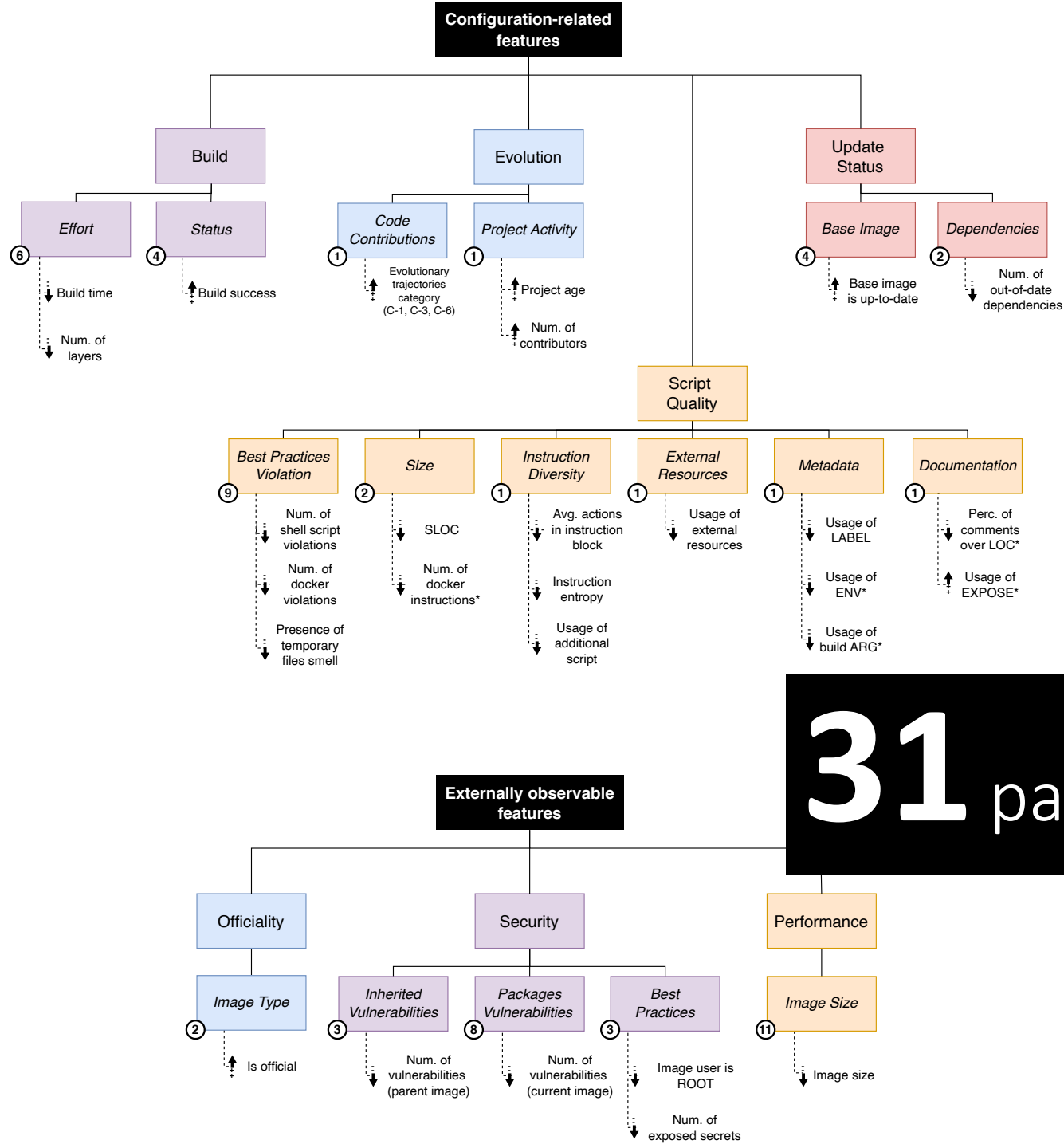"community images are more resource-efficient and have fewer vulnerabilities"

What **quality features** characterize the adoption of a Docker image (and its Dockerfile)

# Step 1:
# Learning from the Literature

**Literature review**

**31** papers

**Configuration-related features**

- **Build**
  - *Effort* (6)
    - ↓ Build time
    - ↓ Num. of layers
  - *Status* (4)
    - ↑ Build success
- **Evolution**
  - *Code Contributions* (1)
    - ↕ Evolutionary trajectories category (C-1, C-3, C-6)
  - *Project Activity* (1)
    - ↑ Project age
    - ↑ Num. of contributors
- **Update Status**
  - *Base Image* (4)
    - ↑ Base image is up-to-date
  - *Dependencies* (2)
    - ↓ Num. of out-of-date dependencies

- **Script Quality**
  - *Best Practices Violation* (9)
    - ↓ Num. of shell script violations
    - ↓ Num. of docker violations
    - ↓ Presence of temporary files smell
  - *Size* (2)
    - ↓ SLOC
    - ↓ Num. of docker instructions*
  - *Instruction Diversity* (1)
    - ↓ Avg. actions in instruction block
    - ↓ Instruction entropy
    - ↓ Usage of additional script
  - *External Resources* (1)
    - ↓ Usage of external resources
  - *Metadata* (1)
    - ↓ Usage of LABEL
    - ↓ Usage of ENV*
    - ↓ Usage of build ARG*
  - *Documentation* (1)
    - ↓ Perc. of comments over LOC*
    - ↕ Usage of EXPOSE*

**Externally observable features**

- **Officiality**
  - *Image Type* (2)
    - ↑ Is official
- **Security**
  - *Inherited Vulnerabilities* (3)
    - ↓ Num. of vulnerabilities (parent image)
  - *Packages Vulnerabilities* (8)
    - ↓ Num. of vulnerabilities (current image)
  - *Best Practices* (3)
    - ↓ Image user is ROOT
    - ↓ Num. of exposed secrets
- **Performance**
  - *Image Size* (11)
    - ↓ Image size

# Externally observable features

## Taxonomy of quality metrics

**6** metrics

*e.g.*

Has "official image" badge

Image size (on disk)

Num. of security vulnerabilities
….

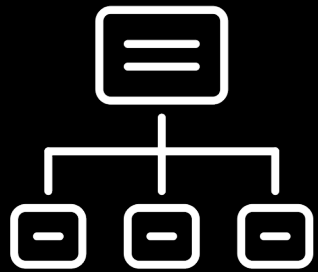# Configuration-related features

**Taxonomy of quality metrics**

**22** metrics

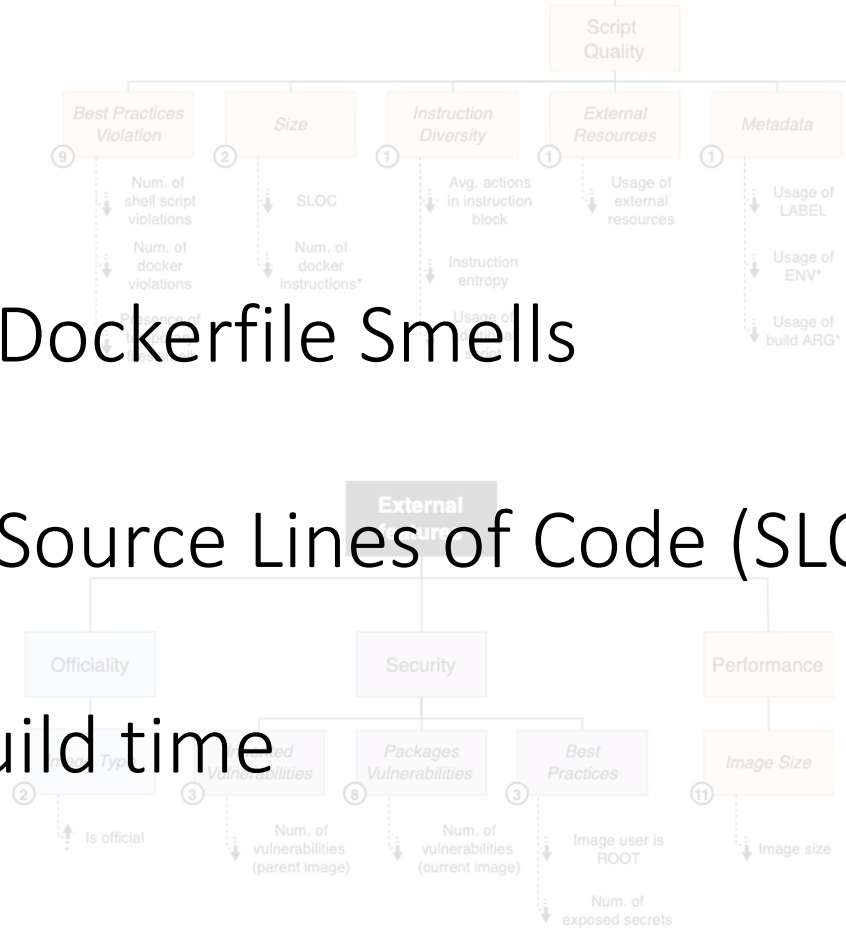*e.g.*

Num. of Dockerfile Smells

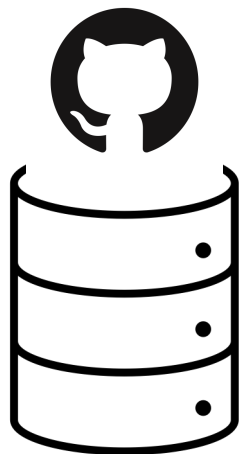Num. of Source Lines of Code (SLOC)

Image Build time

...

# Step 2:
# **Catching the Developers' Preferences**

# RQ1

Can the externally observable features explain the developers' preference for a Docker image?
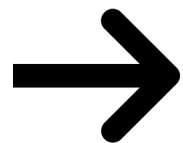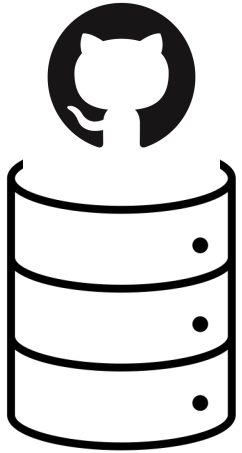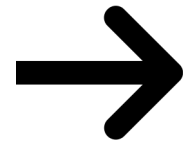
~**2.4k** Docker images

→

~**50k**

open-source
repos

**RQ1: Experiment**

Docker images with
**quality metrics**

$\longrightarrow$

GLM

Docker images
with
**quality metrics**

2 x GLM

R.E. App name
App version

RQ1: Results

Image user is ROOT

Number of vulnerabilities

Image size

Number of exposed secrets

Is official

Number of DockerHub stars

Number of adoptions
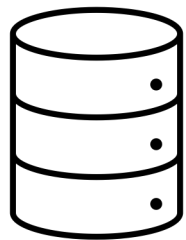
RQ1: Results

# RQ2

Are configuration-related features correlated with externally observable features ?

**2.4k**
Docker images

**10**
most-used apps
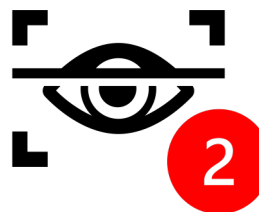
→

**~300** with
source Dockerfiles

RQ2: Results

RQ2: Results

RQ2: Results

**Takeaways**

Developers mostly adopt **official** images

Project age

Image user is ROOT

Usage of additional script
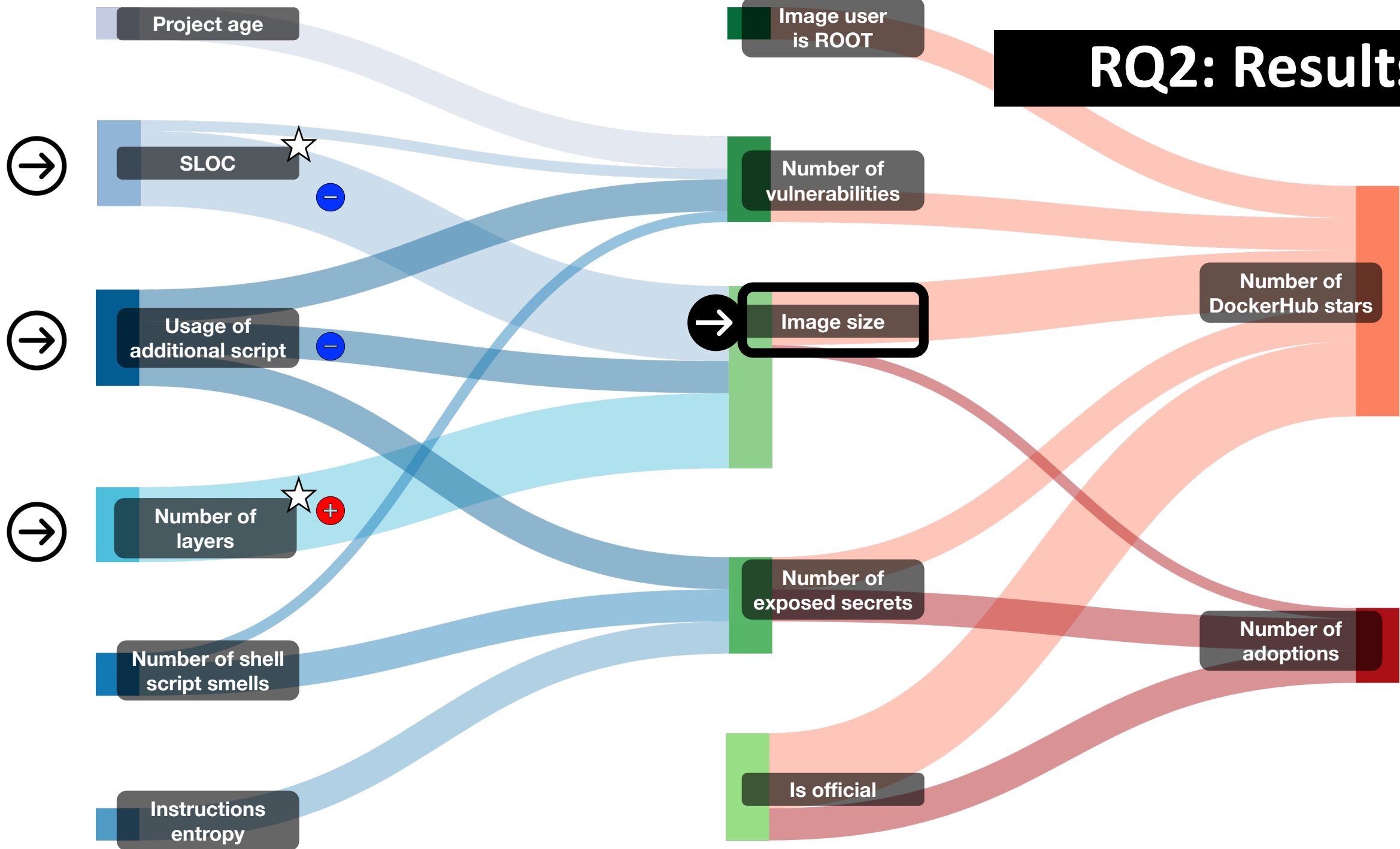
Image size

Number of layers

Number of exposed secrets

Number of shell script smells

Number of DockerHub stars

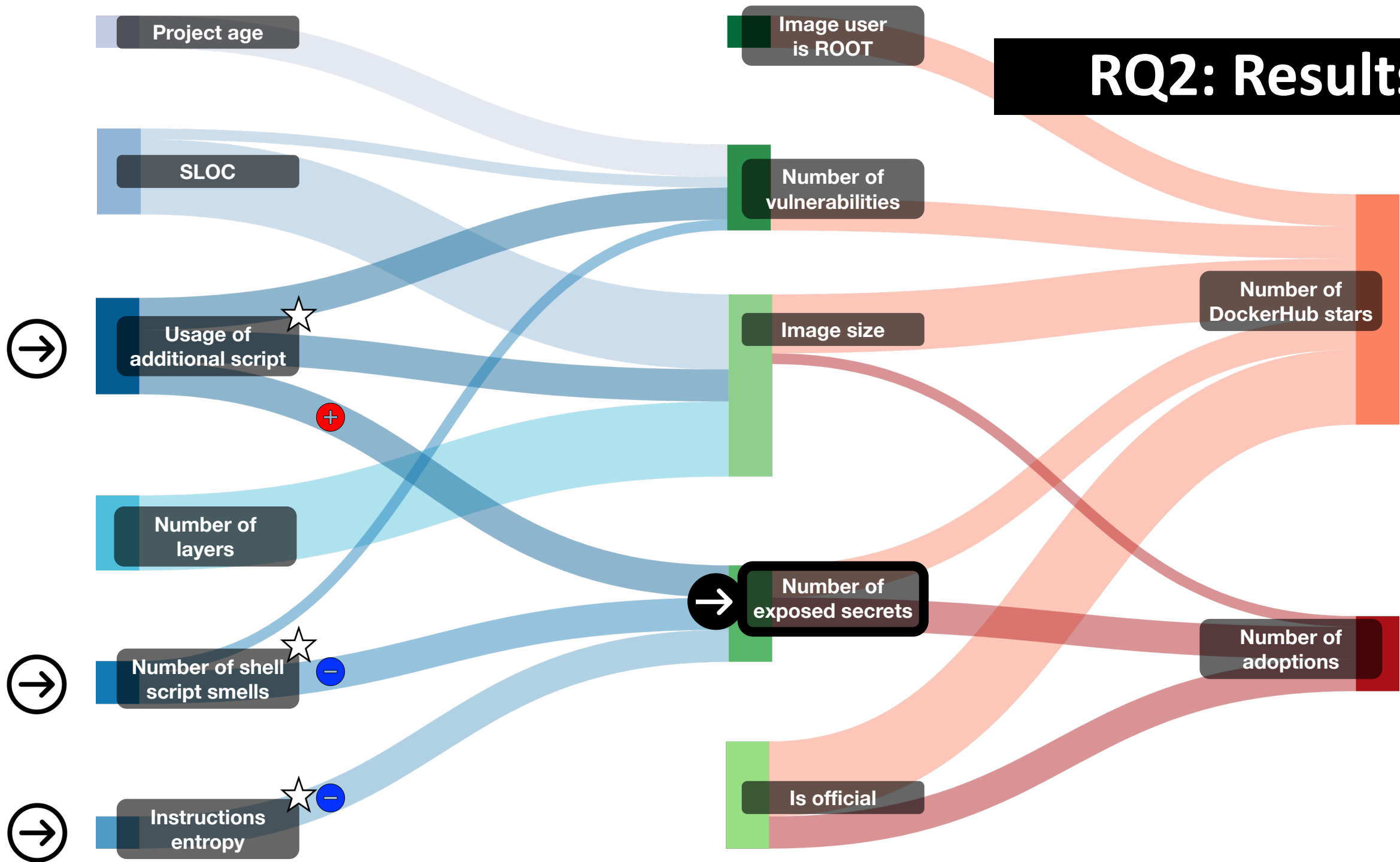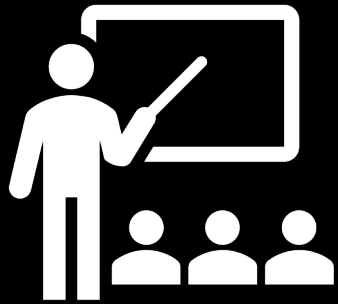Instructions entropy

Is official

Number of adoptions

Takeaways

Developers mostly adopt **official** images

**Fewer SLOC** not means a **lower** image **size**

# Summary



Which Docker image to choose?

## RQ1: Context

~50k open-source repos → ~2.4k Docker images

`1 FROM node:12-alpine`

App name · Version · Flavour

10 most-used apps

## RQ1: Experiment

Docker images with **quality metrics** → 2 x GLM

R.E. App name App version

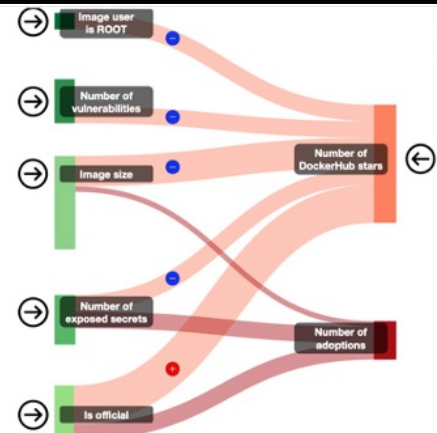"I **like** this image" ①

"I am **using** this image" ②

## RQ1: Results



## RQ2: Experiment

3 x GLM

R.E. App name App version

# vulnerabilities ①
# exposed secrets ②
# image size ③

## RQ2: Results



# Giovanni Rosa

@ **giovanni.rosa@unimol.it**

Get the paper here! →