# Automatically Generating Dockerfiles via Deep-Learning: Challenges and Promises

**Giovanni Rosa**, Antonio Mastropaolo, Simone Scalabrino, Gabriele Bavota and Rocco Oliveto

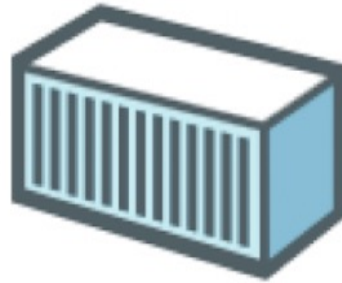**@giovannipink**

**University of Molise, Italy**

17th International Conference on
Software and System Processes
14-15 May 2023, MCEC (Melbourne)

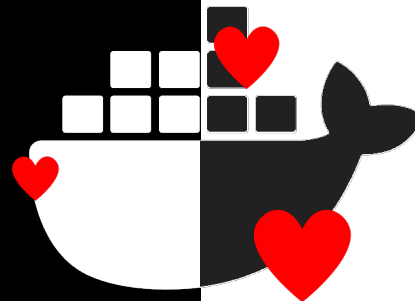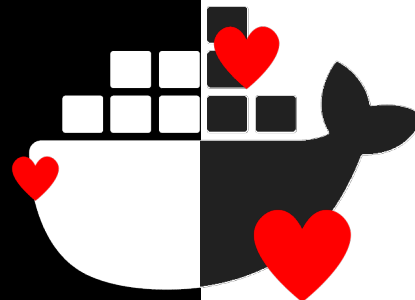Build → Ship → Run

Software Containers

# 85% of

organizations will
adopt containers
by **2025**

Gartner®

# 85% of organizations will adopt containers by 2025

**Gartner®**

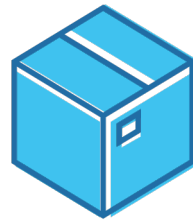# #1 Most-Wanted and #1 Most Loved tool

**2022 Developer Survey**

Docker in a Nutshell

Docker in a Nutshell

Docker in a Nutshell

```
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
```
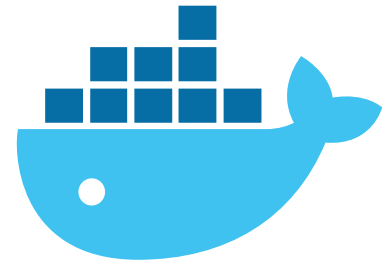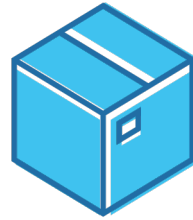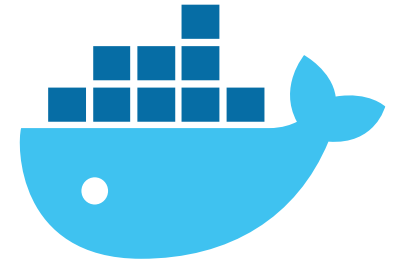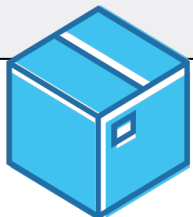
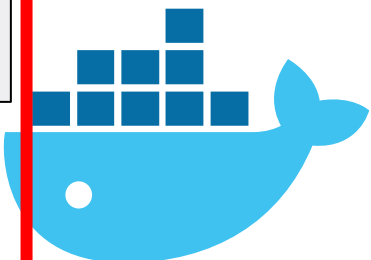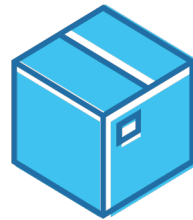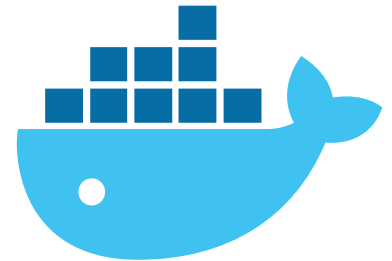**build**

**run**

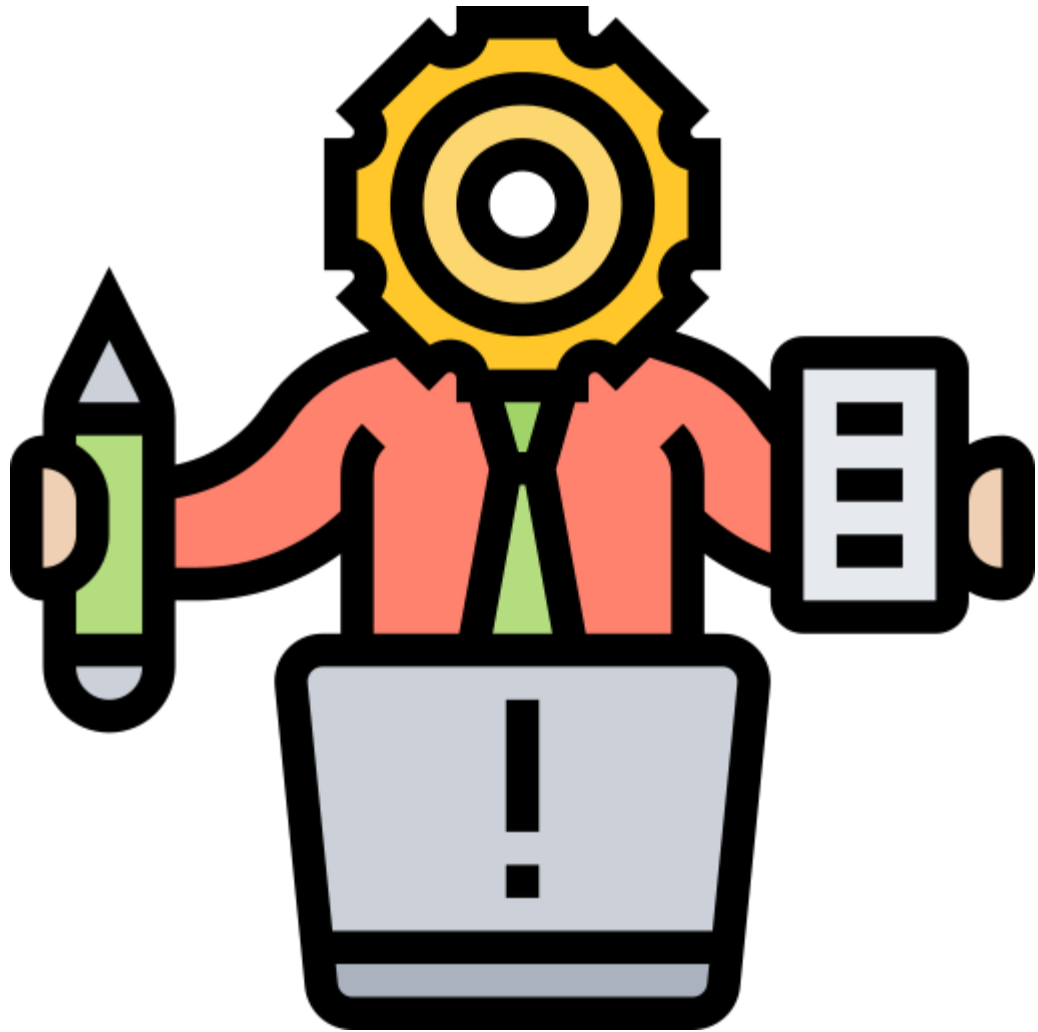Dockerfile

Image

Container

Docker in a Nutshell

Writing
Dockerfiles
is **challenging**

# Time-consuming activity

Reis et. al 2021

Henkel et. al 2020

Time-consuming activity

Technical knowledge required

Reis et. al 2021

**Developing Docker and Docker-Compose Specifications: A Developers' Survey**

Henkel et. al 2020

**Learning from, Understanding, and Supporting DevOps Artifacts for Docker**

Time-consuming activity

Technical knowledge required

**Limited supporting tools**

AI-based
Code Assistants

Horton et. al 2019

Hanayama et. al 2020

Ye et. al 2021

Horton et. al 2019

DockerizeMe: Automatic Inference of Environment Dependencies for Python Code Snippets

Hanayama et. al 2020

Humpback: Code Completion System for Dockerfiles Based on Language Models

Ye et. al 2021

DockerGen: A Knowledge Graph based Approach for Software Containerization

AI-based
**Code Assistants**

...limited to
**dependencies**
and
**code completion**

The automated generation **from scratch** is still an **open challenge**

What about
**ChatGPT?**

**Prompt:**

Can you provide a Dockerfile using **Debian 9** and installing **Golang 1.9.4** via sources?

# Not Bad!

but...

# build error
☹ ☹



```dockerfile
# Use Debian 9 as the base image
FROM debian:9

# Install required dependencies
RUN apt-get update && apt-get install -y \
    wget \
    build-essential \
    git

# Download and extract Go 1.9.4
RUN wget https://dl.google.com/go/go1.9.4.linux-amd64.tar.gz \
    && tar -C /usr/local -xzf go1.9.4.linux-amd64.tar.gz \
    && rm go1.9.4.linux-amd64.tar.gz

# Set Go environment variables
ENV GOROOT=/usr/local/go
ENV GOPATH=/go
ENV PATH=$GOROOT/bin:$GOPATH/bin:$PATH

# Create a workspace directory for Go projects
RUN mkdir -p $GOPATH/src $GOPATH/bin

# Set the working directory
WORKDIR $GOPATH/src

# Define the entry point
CMD ["bash"]
```

Copy code

⟳ Regenerate response

What about
**state-of-the-art DL models**
for code-related tasks?

# Phase 1:
# **Model Construction**

Natural language requirements → *input* → T5 Model → *prediction* → Generated Dockerfile

Dockerfile Generation via T5

**Operating System**: "alpine"

**Package Manager**: "apk"

**Package Requirements**: ["python3"]

**Download from Sources**: FALSE

**ENV variables**: FALSE

**Build ARGs**: FALSE

**LABEL**: TRUE

**EXPOSE for ports**: TRUE

**CMD**: TRUE

**ENTRYPOINT**: FALSE

High-Level Specification

**Operating System**: "alpine"
**Package Manager**: "apk"
**Package Requirements**: ["python3"]
**Download from Sources**: FALSE
**ENV variables**: FALSE
**Build ARGs**: FALSE
**LABEL**: TRUE
**EXPOSE for ports**: TRUE
**CMD**: TRUE
**ENTRYPOINT**: FALSE

High-Level Specification

**Operating System**: "alpine"

**Package Manager**: "apk"

**Package Requirements**: ["python3"]

**Download from Sources**: FALSE

**ENV variables**: FALSE

**Build ARGs**: FALSE

**LABEL**: TRUE

**EXPOSE for ports**: TRUE

**CMD**: TRUE

**ENTRYPOINT**: FALSE

HLS

High-Level Specification

```
1 FROM alpine:3.17
2
3 LABEL maintainer="happy@container.net"
4
5 # Install python3
6 RUN apk add --update --no-cache python3
7
8 COPY . .
9
10 RUN pip install -r requirements.txt
11
12 WORKDIR /tiddlywiki
13
14 EXPOSE 8000
```

**OS**:
"alpine3.17"

**Pkg. Manager**:
"apk"

**Pkg. Requirements**:
"python3"

# High-Level Specification

+2 years

12 software developers

> 50% agrees with the **requirements specification**

**Eng et. al 2021**

# 9.4M
unique Dockerfiles

from **2013** to **2020**

Dockerfile Dataset

Revisiting Dockerfiles in Open Source Software Over Time

Kalvin Eng
Department of Computing Science
University of Alberta
Canada
kalvin.eng@ualberta.ca

Abram Hindle
Department of Computing Science
University of Alberta
Canada
abram.hindle@ualberta.ca

Abstract—Docker is becoming ubiquitous with containerization for developing and deploying applications. Previous studies have analyzed Dockerfiles that are used to create container images in order to better understand how to improve Docker tooling. These studies obtain Dockerfiles using either Docker Hub or Github. In this paper, we revisit the findings of previous studies using the largest set of Dockerfiles known to date with over 9.4 million unique Dockerfiles found in the World of Code infrastructure spanning from 2013-2020. We contribute a historical view of the Dockerfile format by analyzing the Docker engine changelogs and use the history to enhance our analysis of Dockerfiles. We also reconfirm previous findings of a downward trend in using OS images and an upward trend of using language images. As well, we reconfirm that Dockerfile smell counts are slightly decreasing meaning that Dockerfile authors are likely getting better at following best practices. Based on these findings, it indicates that previous analyses from prior works have been correct in many of their findings and their suggestions to build better tools for Docker image creation are further substantiated.

Index Terms—Git, GitHub, Docker

## I. INTRODUCTION

Docker, a tool for creating and running programs in containers consistently across platforms, was initially released to the public on March 20, 2013 [1], [2]. Ever since its release, Docker has amassed a considerable following with 2.9 million desktop installations and 7 million Docker Hub users as reported in July 2020 [3].

The use of container software such as Docker has made applications easier to deploy, scale, and migrate across platforms. Furthermore, it has also made development setup simpler by reducing the amount of time needed to configure an appropriate environment by bundling the needed configuration instructions in a *Dockerfile* which can then be used to create images for containers.
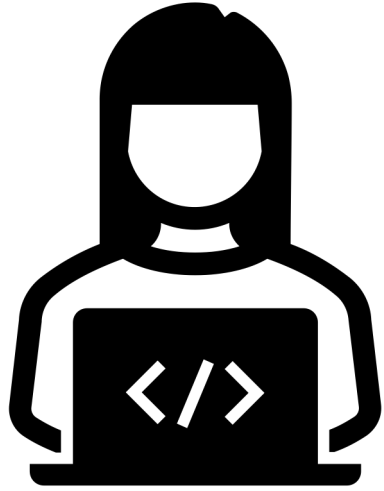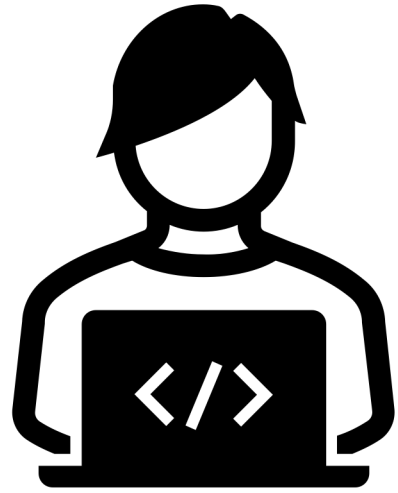
Because of the proliferation of Docker, this paper seeks to replicate and elaborate on previous studies on Dockerfile usage using the largest Dockerfile dataset [4] known to date. This paper has findings, using data between 2013-2020, that include:

- Discovering that 7.99% of Dockerfiles exist in more than one distinct repository
- Most repositories overall contain up to 6 Dockerfiles
- Confirmation of previous findings such as JavaScript being the most popular language of projects that contain

Dockerfiles [5], [6] (2016, 2020) and RUN being the most popular Dockerfile instruction [5].

## II. PREVIOUS WORK

In previous work, large collections of Dockerfiles have been mined from Github and Docker Hub to better understand Docker use in repositories and to gather insights on popularity, quality, and possible ways to improve Docker usage.

*Mining Github:* Cito et al. [5] (2016) focused on analyzing over 70,000 Dockerfiles in Github within commits up until October 2016 finding that: most Dockerfiles use heavy-weight operating systems as a base image; the biggest quality issue of Dockerfiles is missing version pinning of images; and Dockerfiles are not revised often. In another study by Wu et al. [7] (2020), 6334 projects were selected from Github and analyzed for Dockerfile smells finding that: 62% of projects selected have code smells; newer and popular projects have less code smells; and projects with different languages have discernible differences in the amount of smells. Also of note is Henkel et al. [8] who retrieved approximately 178,000 Dockerfiles from Github to test with rules mined from the Dockerfiles of official Docker images and found that there should be more tooling to support developers using Dockerfiles.

*Mining Docker Hub:* Lin et al. [6] (2020) scraped Docker Hub and its related GitHub and Bitbucket repositories retrieving 434,304 Dockerfiles up until May 2020. They sought to better understand the Docker ecosystem through Docker Hub. They concluded that: for base images more programming runtime images and ready-to-use application images are being used instead of OS images; there is a declining trend over the years in Dockerfile smells; and there is an upward trend of using end of life Ubuntu base images. Additionally, Zhang et al. [9], [10] selected 2840 projects from Docker Hub to identify evolutionary patterns of Dockerfiles and its impact on Dockerfile quality and image build latency. It should be noted that mining from Docker Hub may not be representative of all Docker usage as users do not have to push images to Docker Hub to use Docker and can choose to build and host images locally or in a private repository.
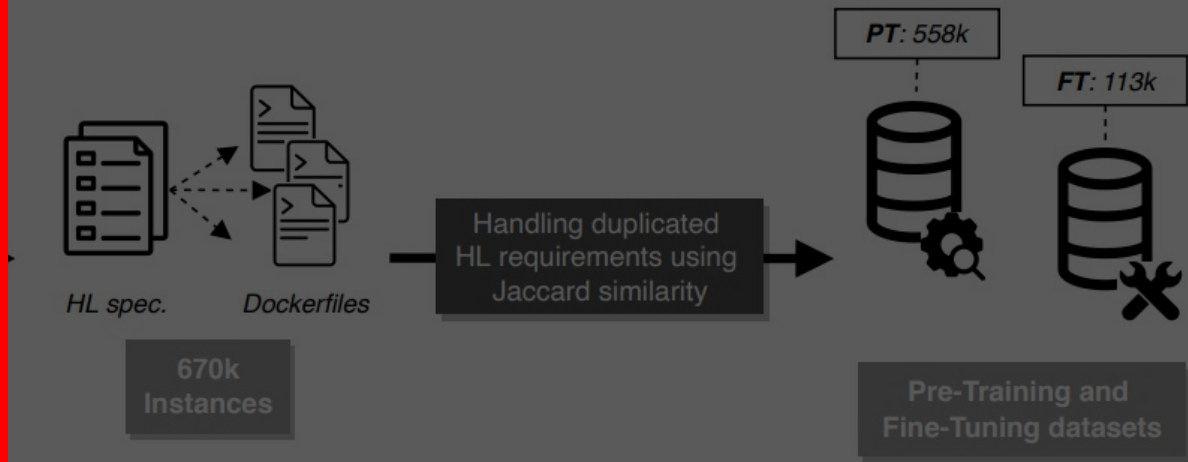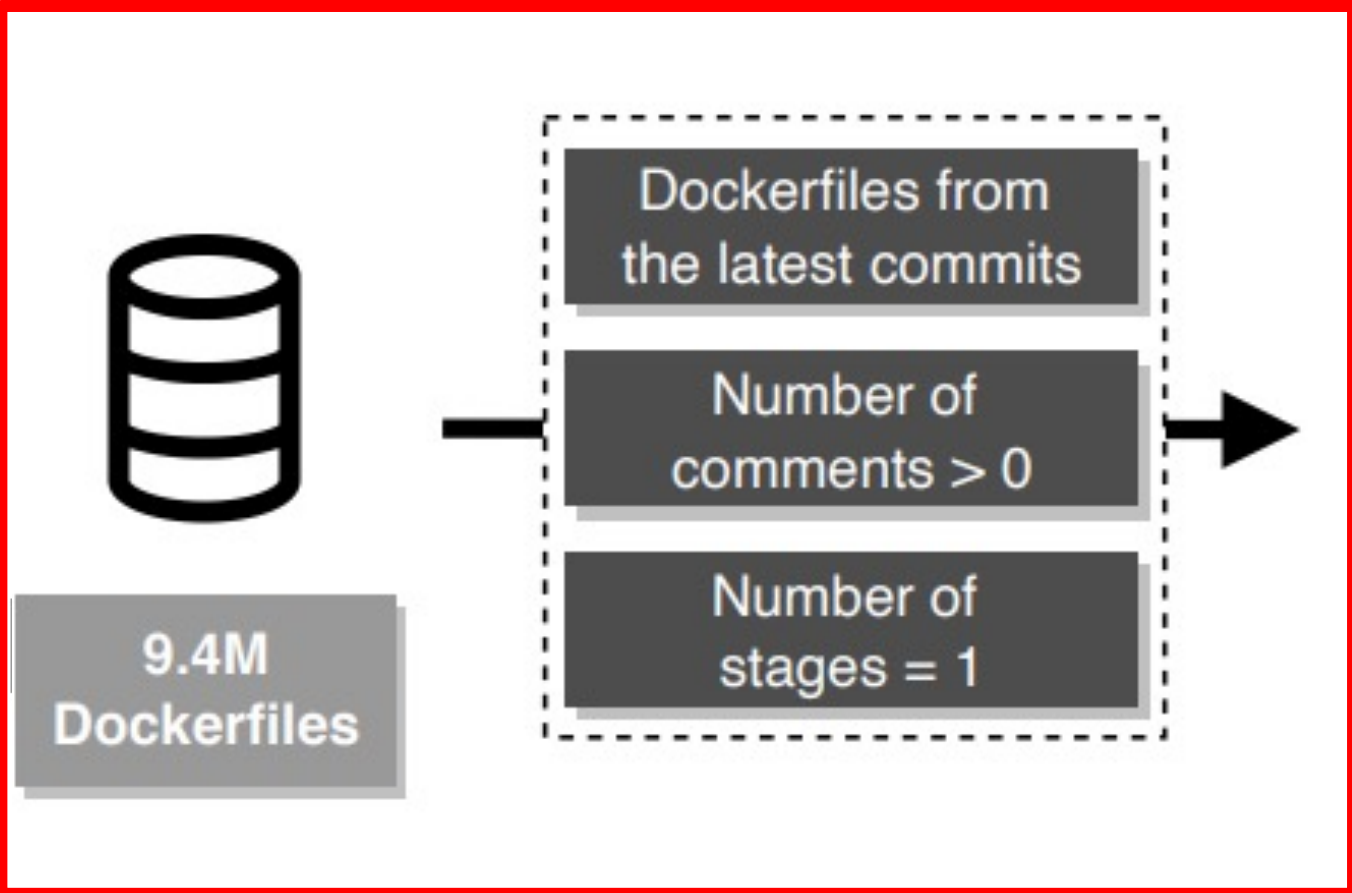
*A. Challenges in Previous Work*

All of the above previous work focuses on Docker use in a project based perspective and involves mining Dockerfiles
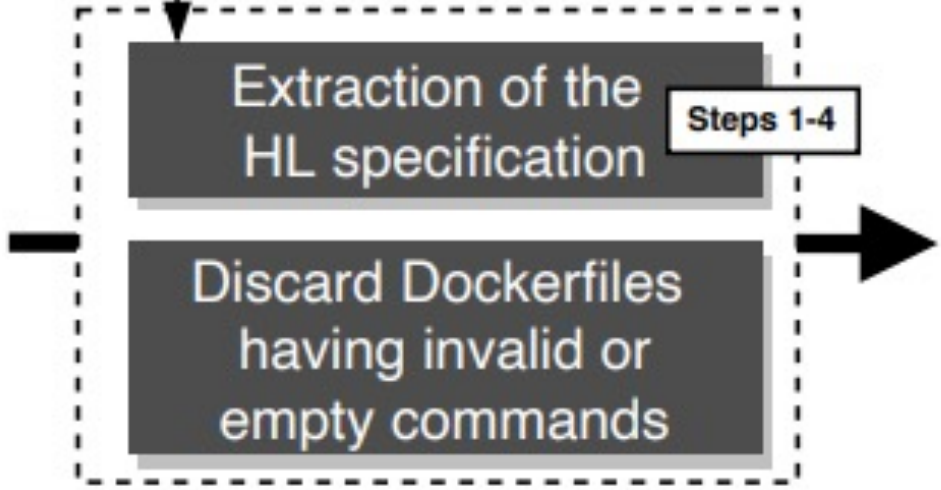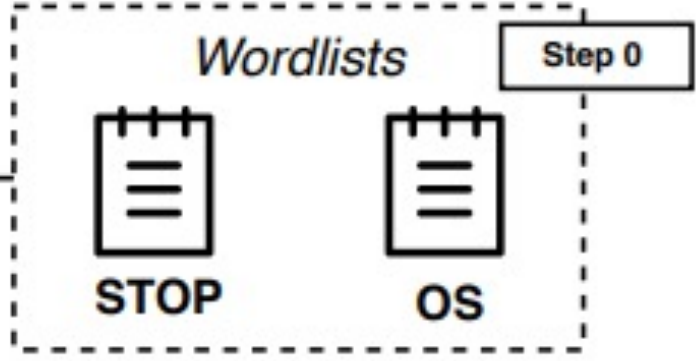
Dataset Filtering

670k Instances

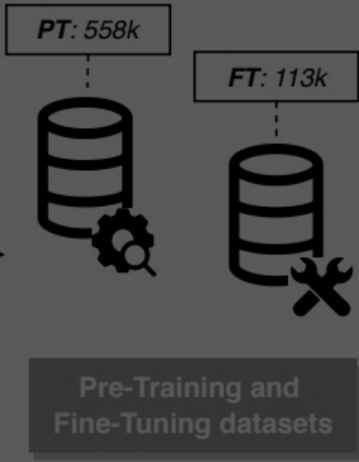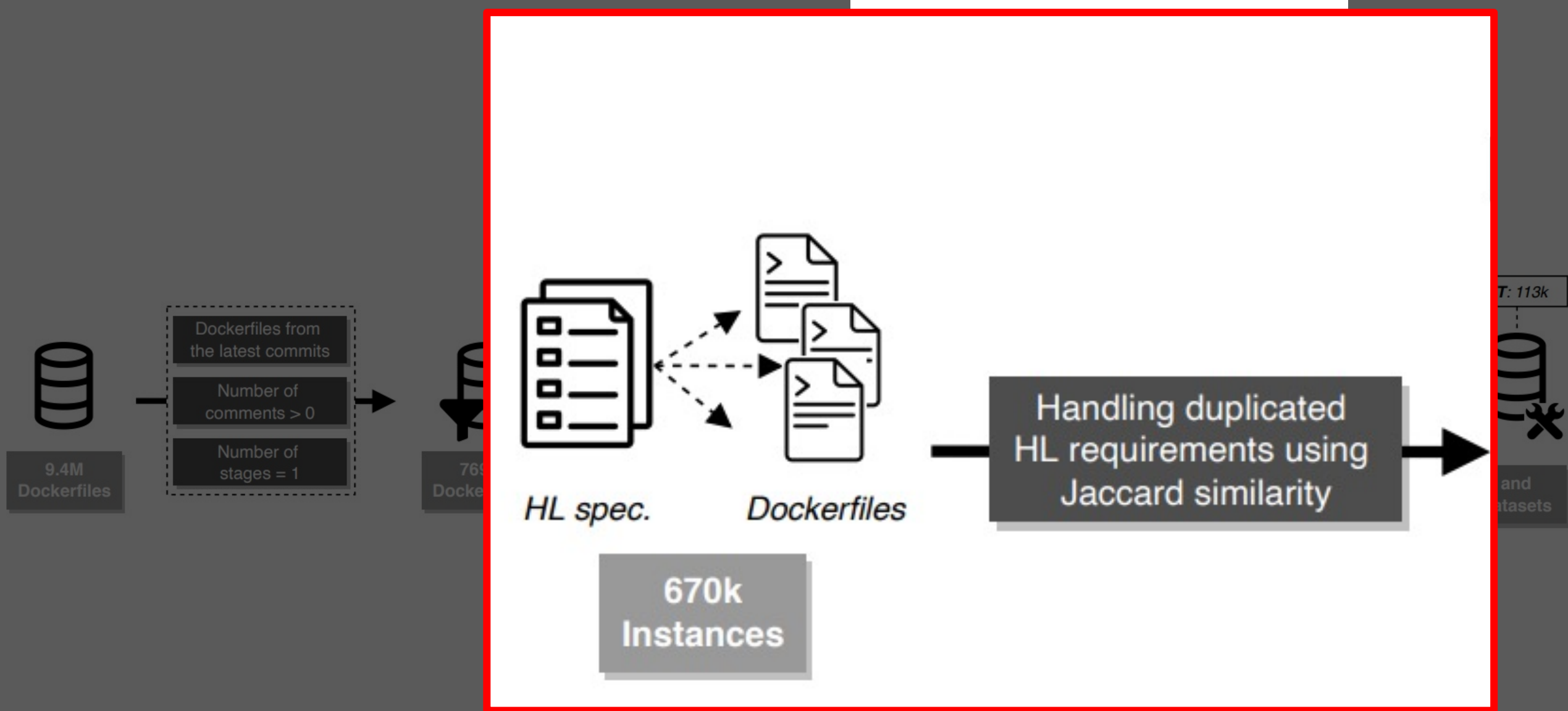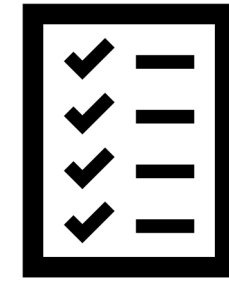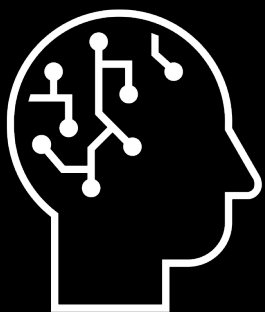HL spec.   Dockerfiles   Handling duplicated HL requirements using Jaccard similarity
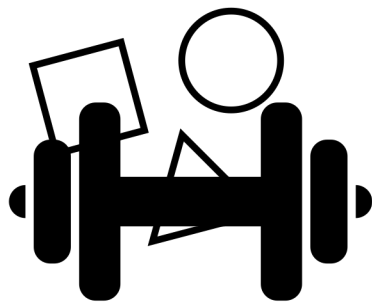
Dataset Filtering

**100k** Model Tuning

**11k** Test
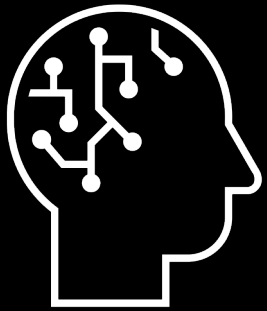
# T5 model construction

Pre-Training

**560k** instances

Fine-Tuning

**90k** instances

# Phase 2:
# **Model evaluation**

# 2 baselines
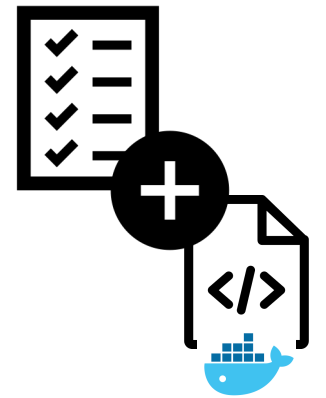
# 90k instances

Elasticsearch

HLS → ES node — bool query →

IR-Baseline 1

SentTransformers

HLS

BERT model

cosine similarity

IR-Baseline 1

**3** dimensions

**11k** Test

|  | T5 | ES | SBERT.net |
|---|---|---|---|
| OS | **~1.00** | 0.92 | 0.88 |
| Pkg. Manager | 0.98 | **1.00** | **1.00** |
| Pkg. Requirements | 0.87 | **0.88** | 0.76 |
| Download from sources | 0.82 | **0.84** | 0.52 |
| ENV variables | **0.89** | 0.81 | 0.17 |
| Build ARGs | **0.99** | 0.88 | 0.17 |
| LABEL | **~1.00** | 0.87 | 0.37 |
| EXPOSE for ports | 0.80 | **0.83** | 0.45 |
| CMD | 0.74 | **0.83** | 0.26 |
| ENTRYPOINT | 0.84 | **0.85** | 0.45 |

# Results for RQ1

```
FROM ubuntu:latest

RUN apt-get update && \
    apt-get install -qqy ...

RUN ./scripts/custom.sh
```

Input Dockerfile

Parsed AST
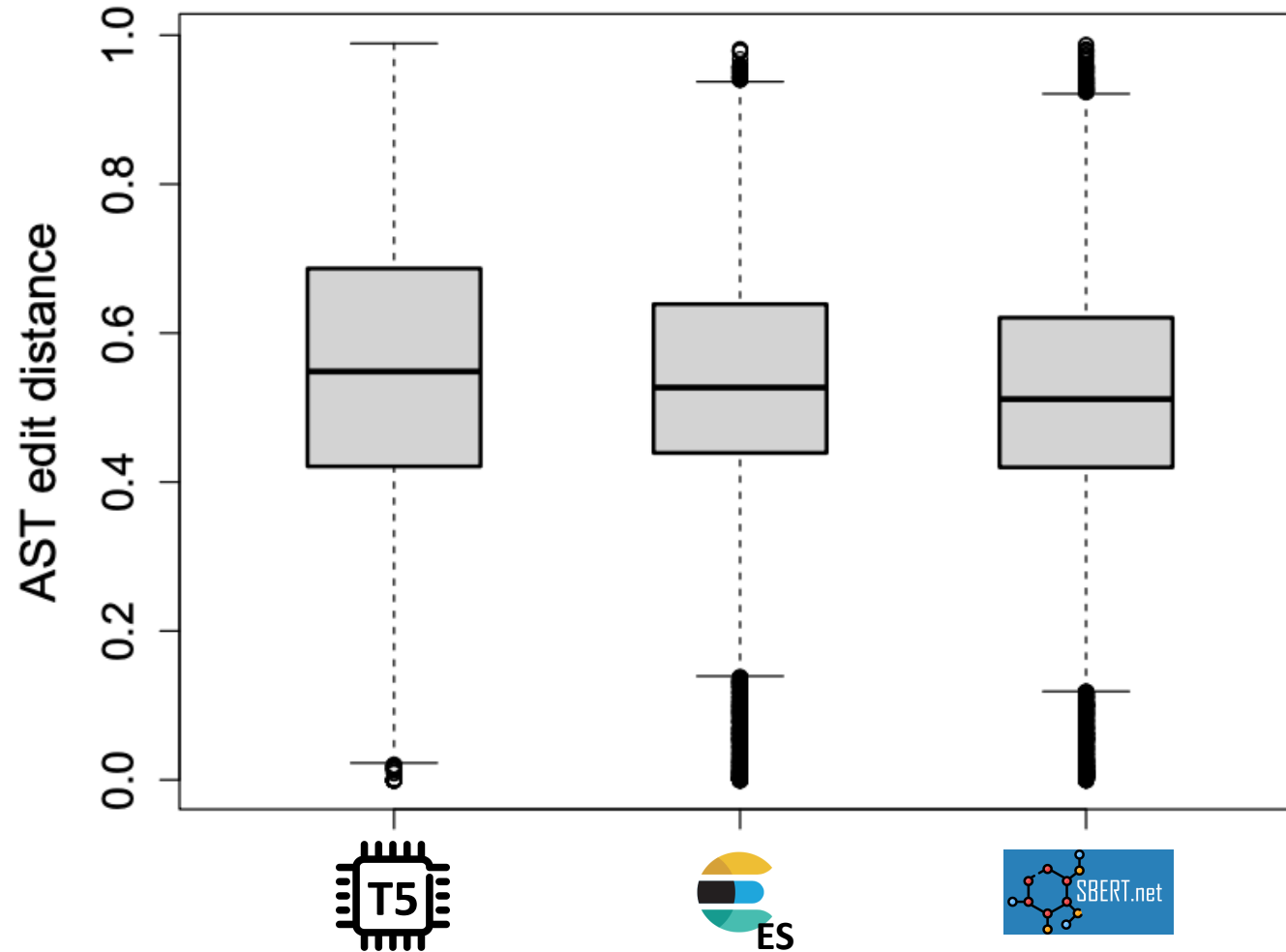(Henkel et. al 2020)

Metric:
**AST
edit distance**

# Summary

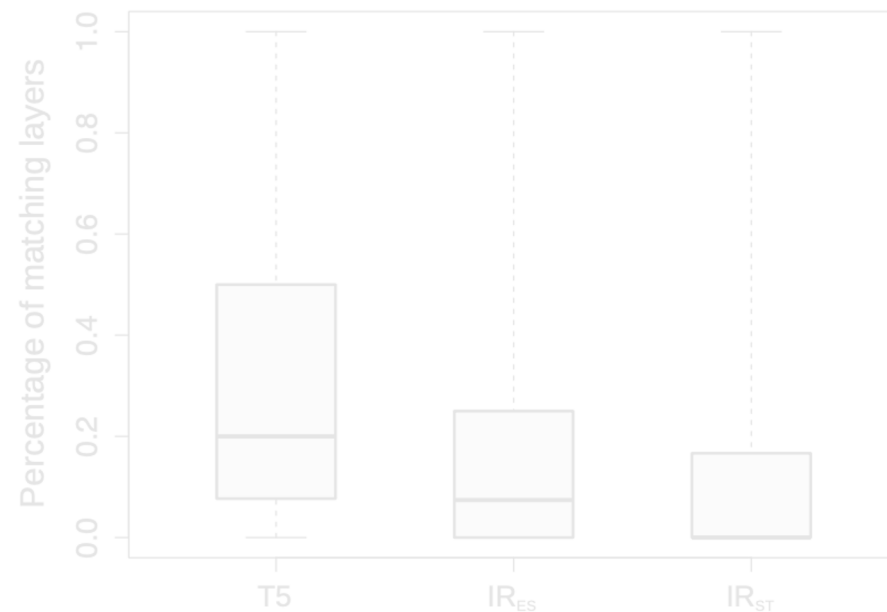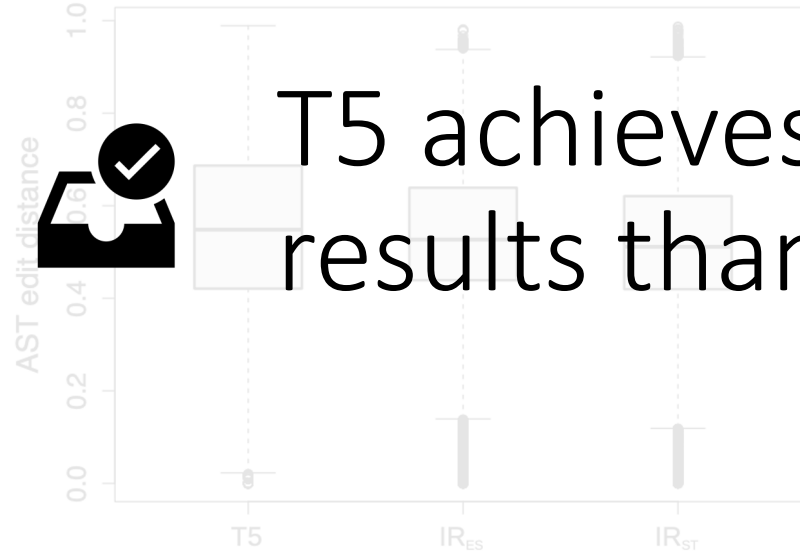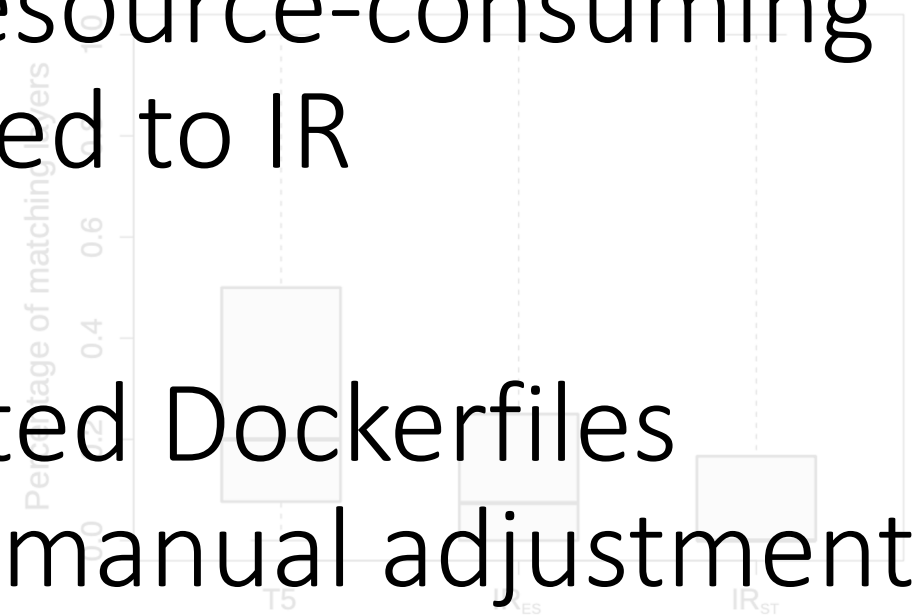T5 achieves slightly better results than IR …

Summary

T5 achieves slightly better results than IR ...
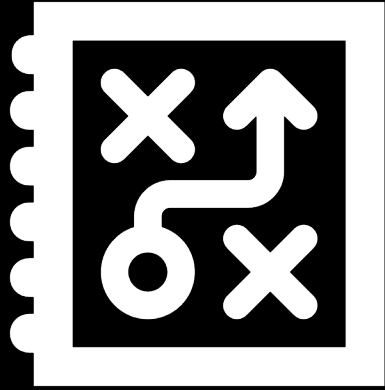
More resource-consuming compared to IR

Generated Dockerfiles require manual adjustments

# What we have learned?

# Challenge
# #1

Not enough training instances

Challenge #2

A different training procedure must be used

# Summary



**Writing** Dockerfiles is **challenging**

Natural language requirements → *input* → T5 Model → *prediction* → Generated Dockerfile

Dockerfile Generation via T5

**3** dimensions

**11k** test

Evaluation

T5 model construction

Fine Tuning

Aa — English Only

🐳 — Dockerfile Only

Aa 🐳 — Dockerfile & English

Summary

T5 achieves slightly better results than IR ...

More resource-consuming compared to IR

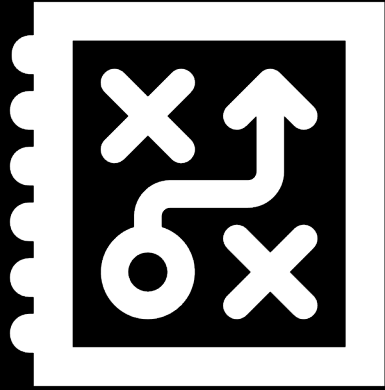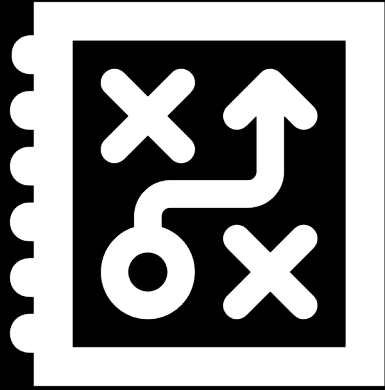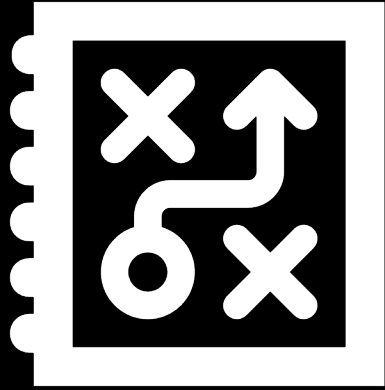Generated Dockerfiles require manual adjustments

**What we have learned?**

**Giovanni Rosa**
https://giovannirosa.com

**Thank you!**