# A Framework for Scientific Workflow Reproducibility in the Cloud

## Rawaa Qasha, Jacek Cała, Paul Watson

Newcastle University, Newcastle upon Tyne, UK

Email: {r.qasha, jacek.cala, paul.watson}@newcastle.ac.uk
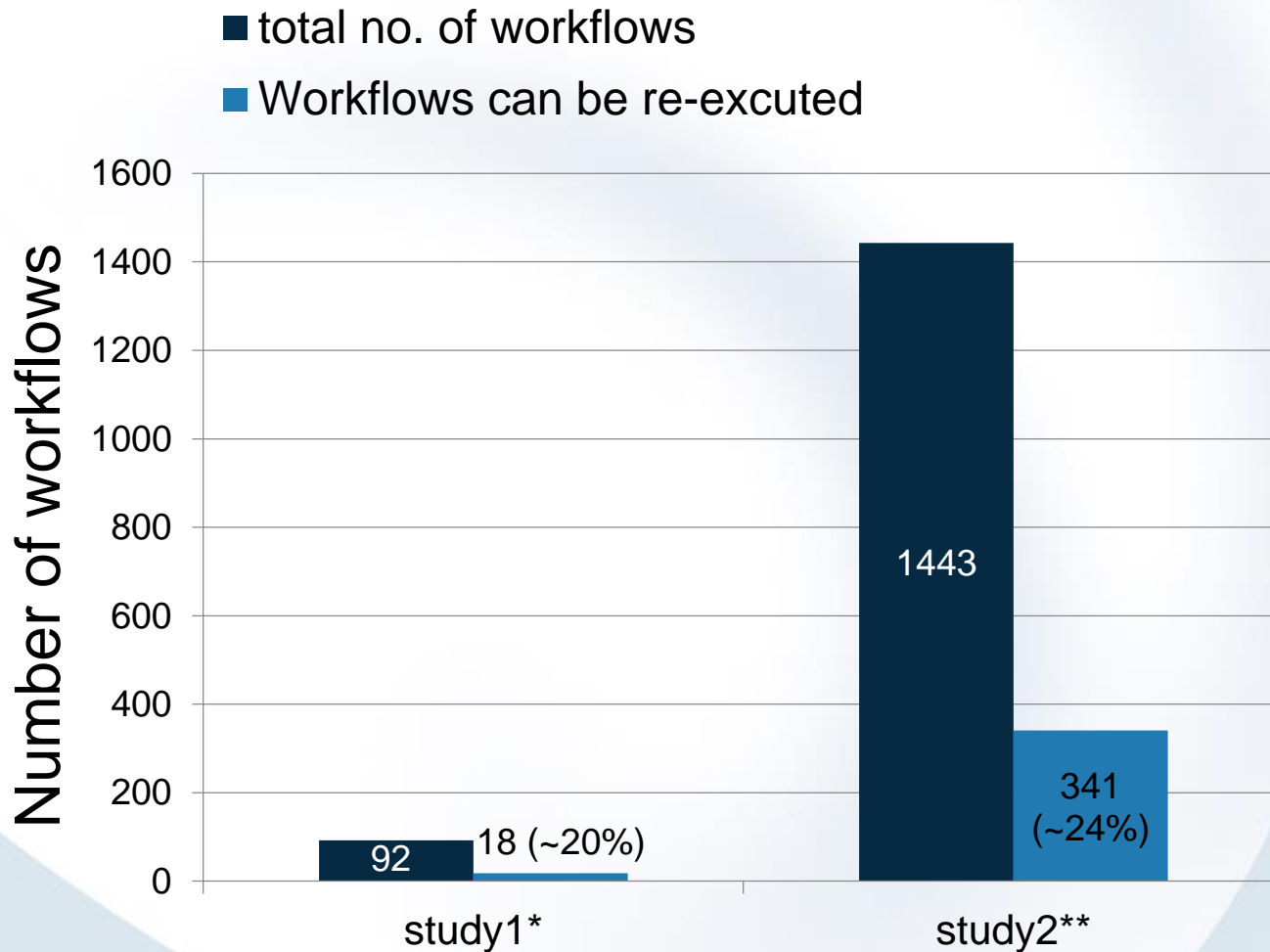
Newcastle University

# In this paper

- A new framework for repeatability and reproducibility of scientific workflow

- Integrating logical and physical preservation approaches

- Offering Workflow/tasks repositories with version control

- Supporting automatic deployment and image capture of workflows and tasks

# Outline

- Background

- Challenges for workflow reproducibility

- Our solution for logical and physical preservations

- Overview of reproducibility framework

- Experiments and results

- Conclusions

# Workflows & Reproducibility



■ total no. of workflows
■ Workflows can be re-excuted

study1* — total: 92, re-executed: 18 (~20%)
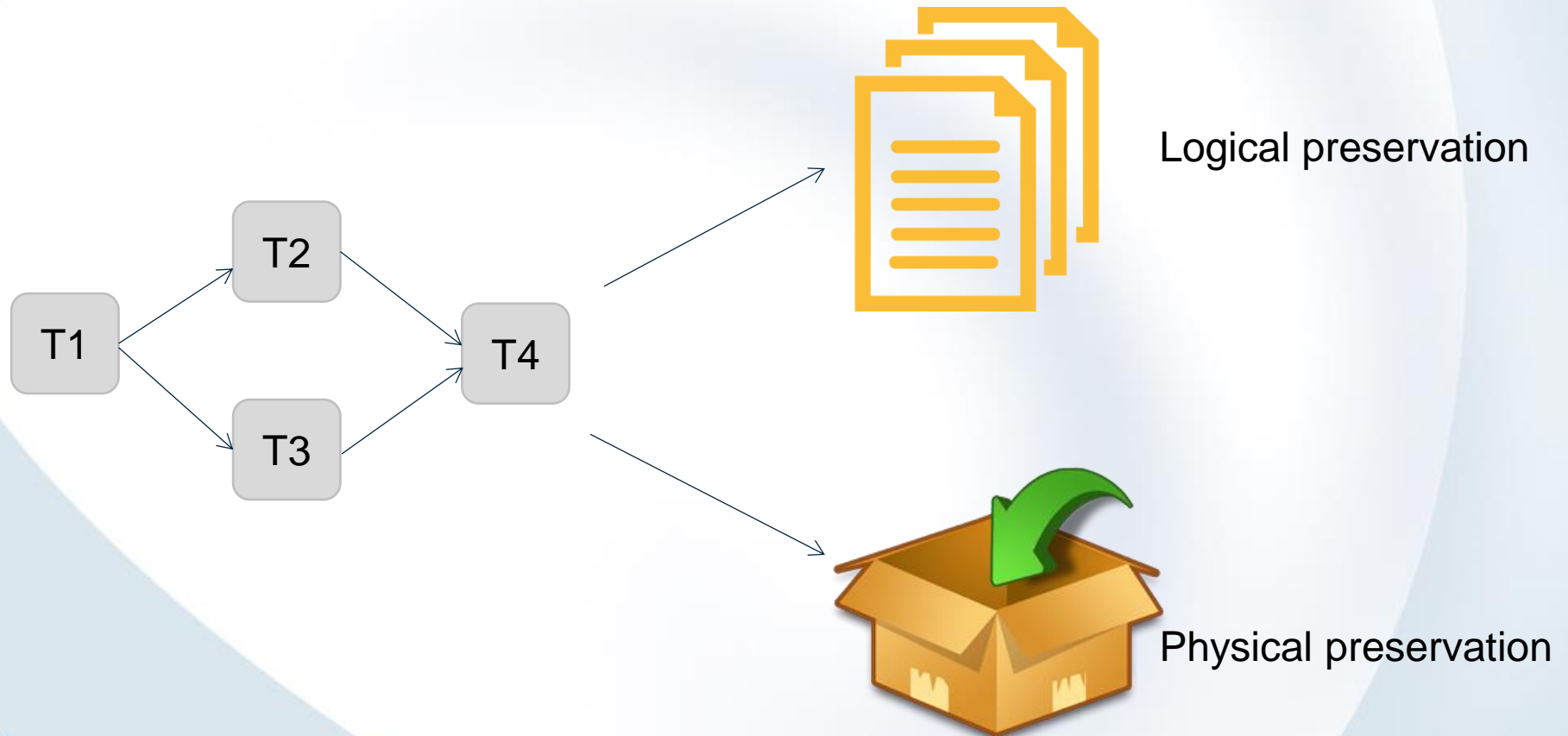study2** — total: 1443, re-executed: 341 (~24%)

*Zhao et al, "Why workflows break Understanding and combating decay in Taverna workflows," 2012
**Mayer et al, "A Quantitative Study on the Re-executability of Publicly Shared Scientific Workflows", 2015
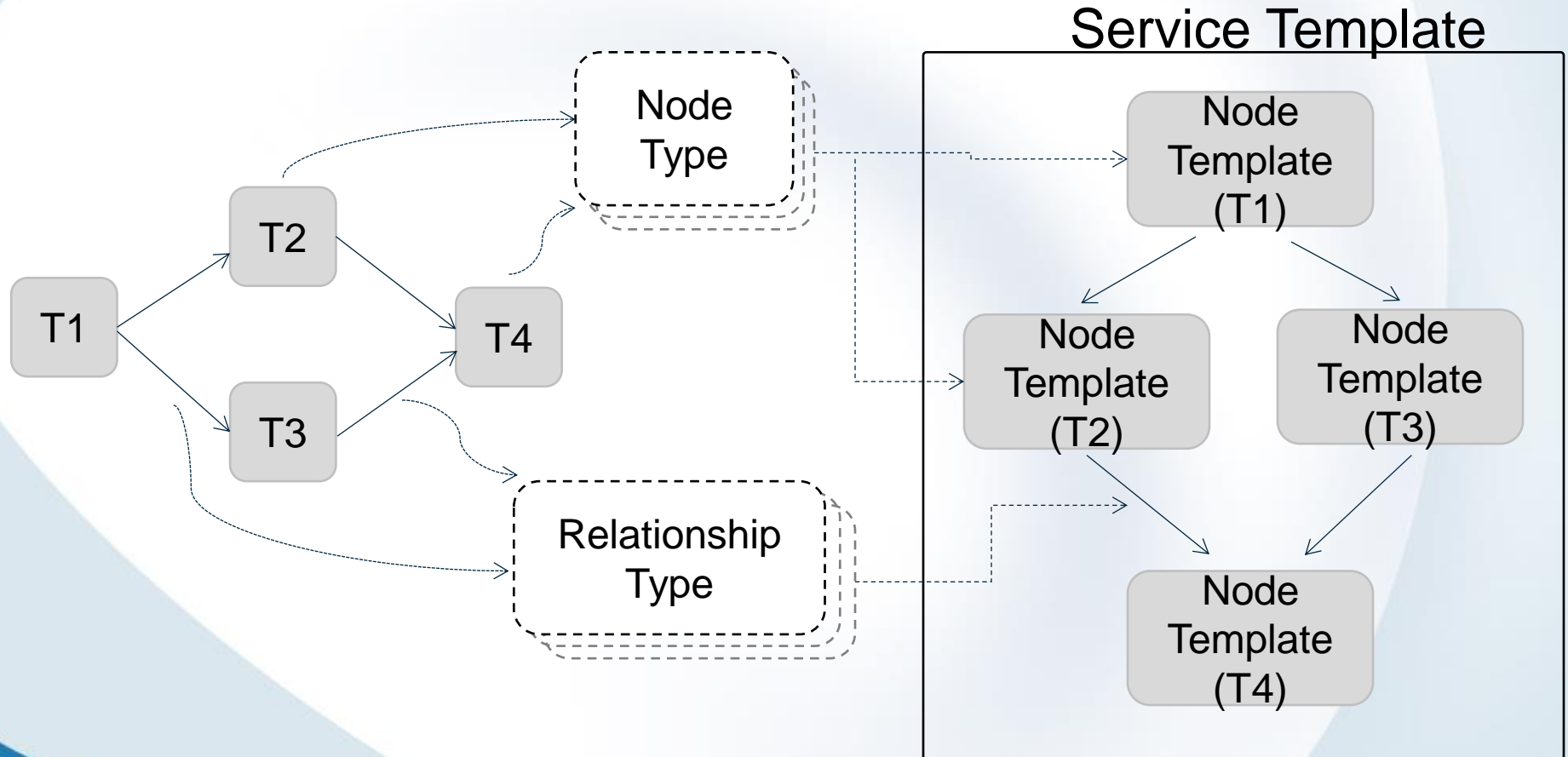
4

# Challenges
# for workflow reproducibility

- Insufficiently detailed workflow description

- Insufficient description of the execution environment

- Unavailable execution environments

- Absence of & changes in the external dependencies
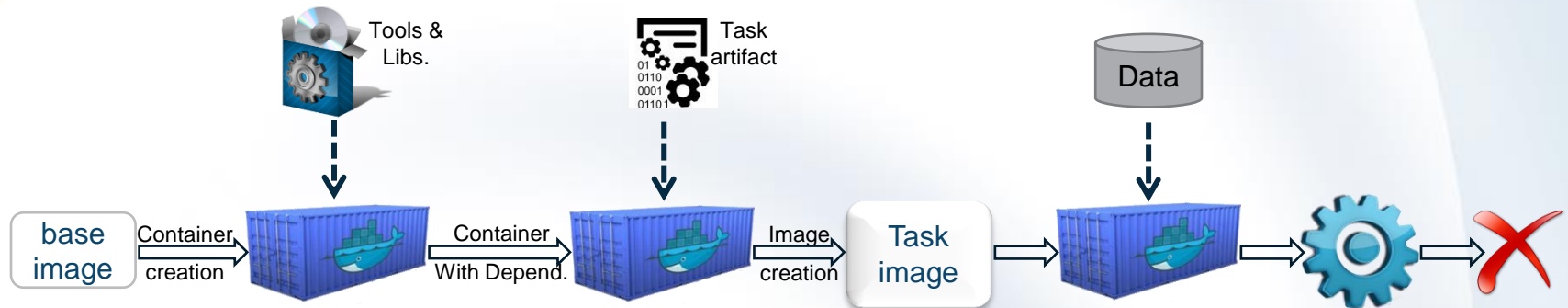
- Missing input data
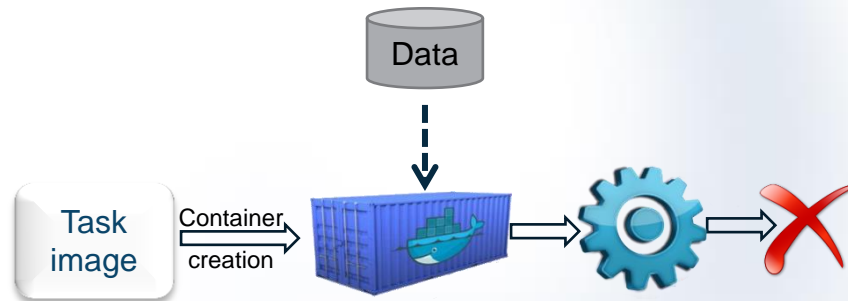
# Common reproducibility approaches

T1

T2

T3

T4

Logical preservation

Physical preservation

# Using TOSCA as a logical preservation

*Workflow and execution environment description*
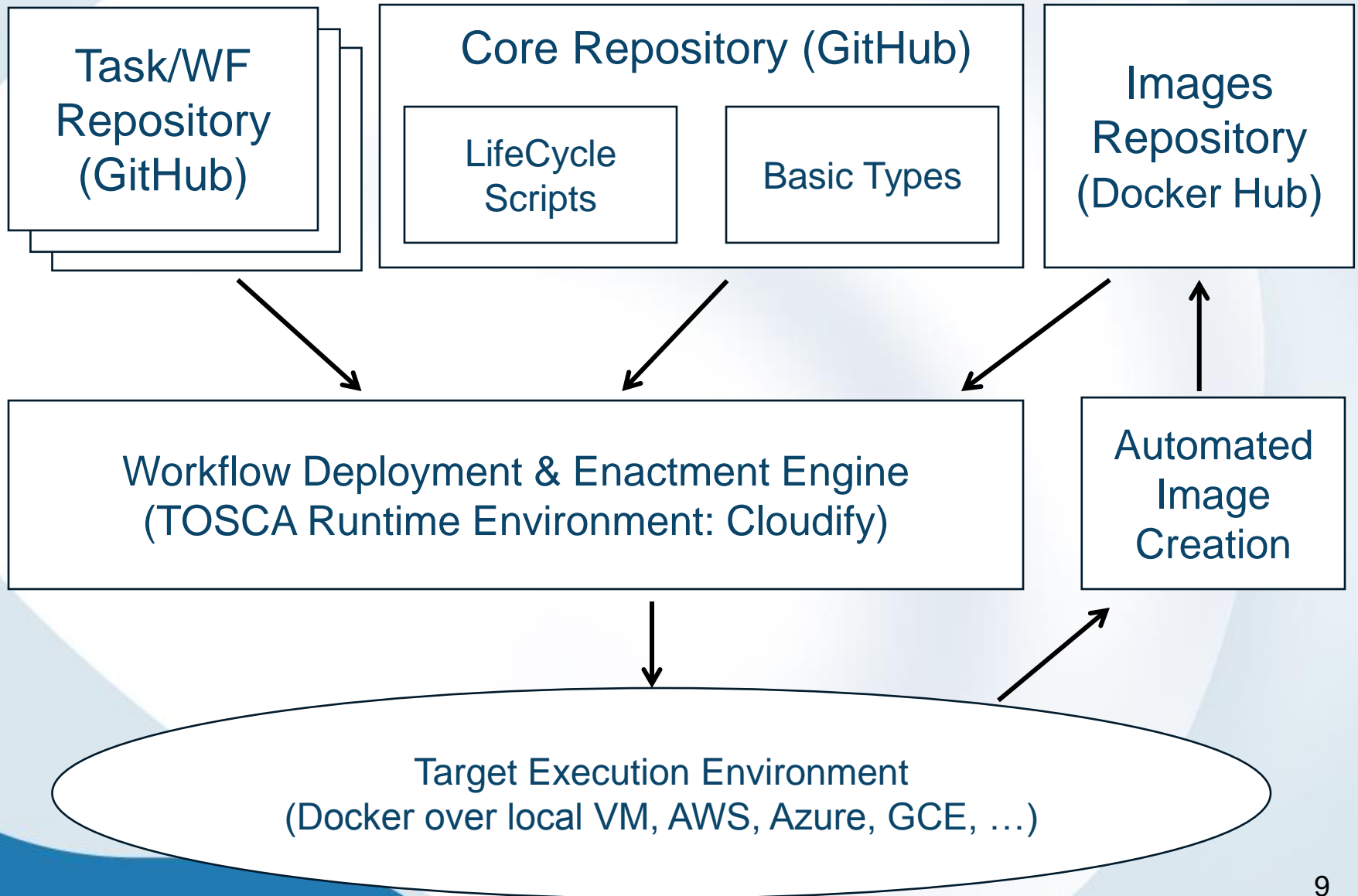
# Using Docker for physical preservation



(a) Initial task deployment & execution



(b) Task deployment & execution with task image

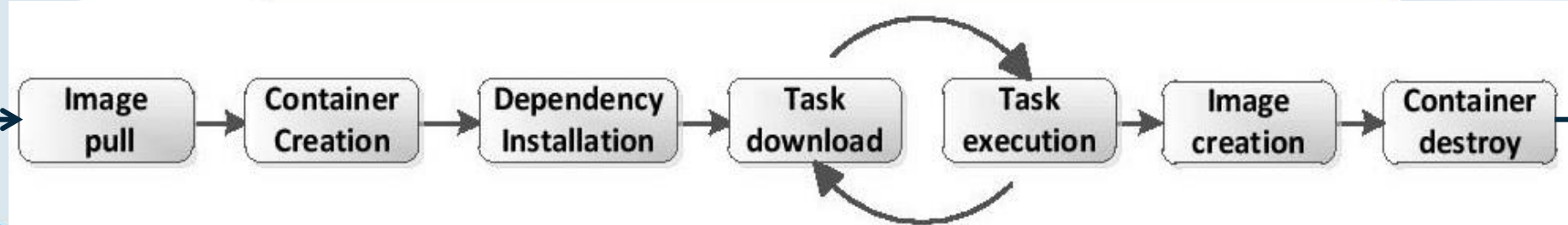*Preserving execution environment and dependencies, tracking changes*

# Reproducibility Framework

Task/WF Repository (GitHub)

Core Repository (GitHub)

LifeCycle Scripts

Basic Types

Images Repository (Docker Hub)

Workflow Deployment & Enactment Engine (TOSCA Runtime Environment: Cloudify)

Automated Image Creation

Target Execution Environment (Docker over local VM, AWS, Azure, GCE, …)

# Multi-container deployment

# Single container deployment

# Time line of workflow devOps

# Workflow repository



**Outputs:**

output-folder: '~/blueprint-name'
output-file(s): {index-BAI-files, output-SAM_BAM-files}
description:
types: {' ', ' '}

**Execution-Environment:**

Cloudify-version: 3.2
Docker-version: 1.8+
OS-type: ubuntu14.04
Disk-space: 10 GB
RAM: 3 GB

## Deployment Instruction

This repository includes all files and scripts to deploy Picard workflow on Multiple Docker containers as follow:

1- Clone the repository to your machine, open a terminal window and change to workflow repository.
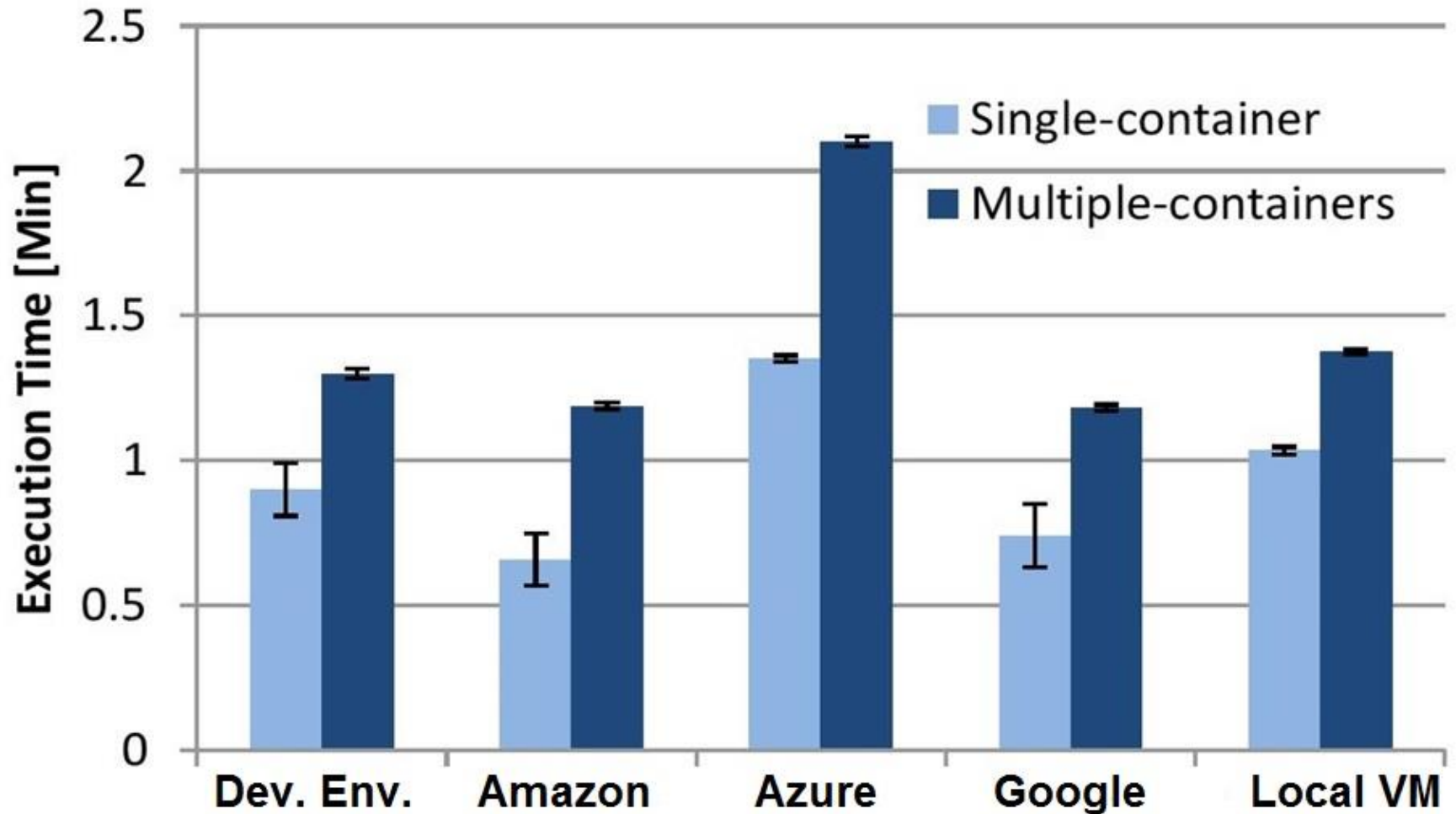2- To execute the workflow with multi containers and the attached input sample, in the terminal run:
. ./Picard-deploy.sh 1
3- If you have own input files, copy your files Dir to Picard/Input-sample folder, open Input.yaml file and change input Dir name, then
run: . ./Picard-deploy.sh 1

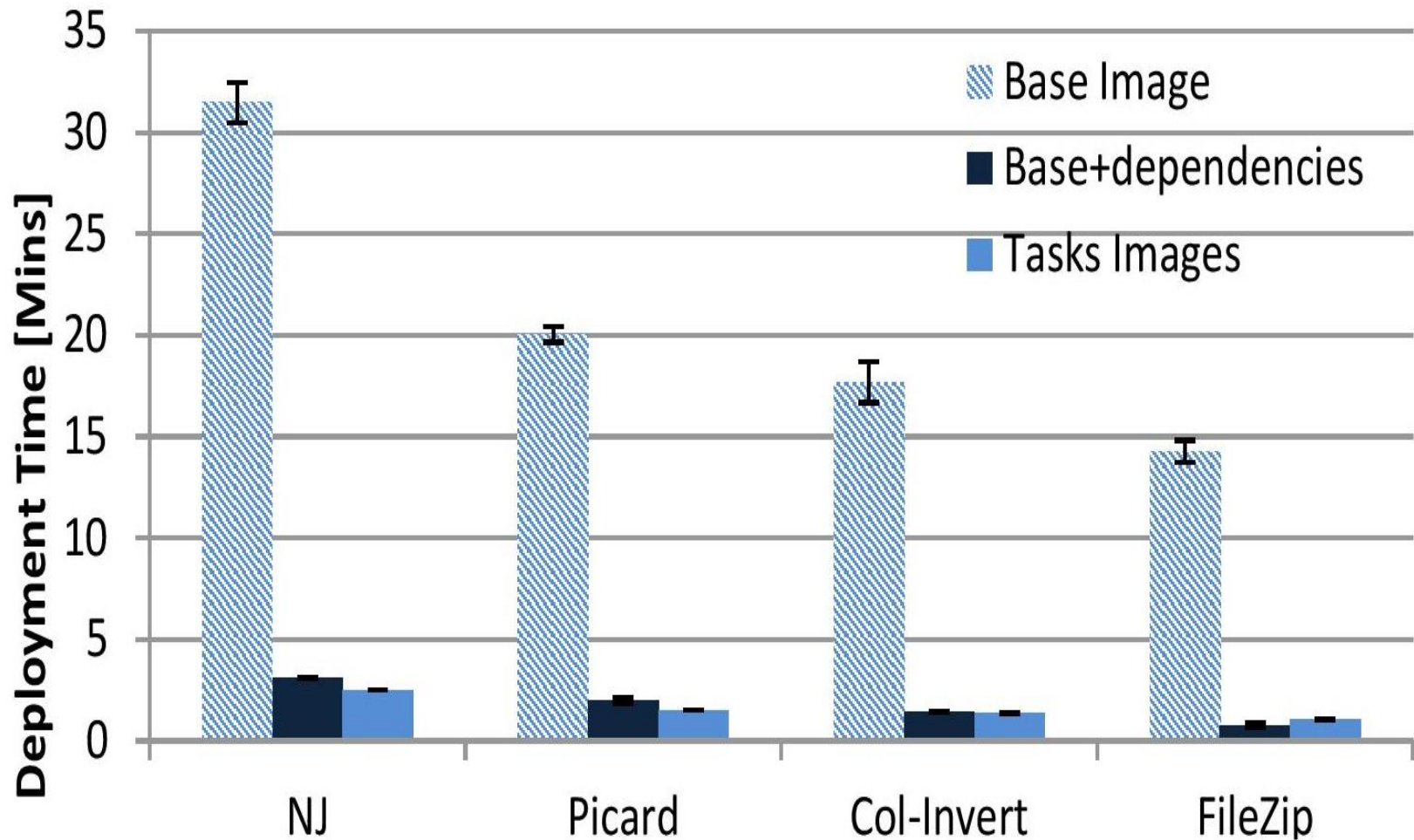4- To execute the workflow with single container, follow either step 2 or 3 but run:
. ./Picard-deploy 2

**Repository file list:**
- 70 commits
- Branch: master
- New pull request
- rawaqasha getting task ID
- Core-LifecycleScripts @ fda31d3
- Input-sample/Data
- scripts
- .gitmodules
- Picard-1container.yaml
- Picard-deploy.sh
- Picard.yaml
- README.md
- input.yaml
- picard.jpg
- picard.png

*Preserving description, input data, tracking changes and deployment instructions*

# Experiments and Results

# 1- Repeatability of a workflow on different clouds

# 2- Automatic image capture for improved performance

# Conclusions

- Full workflow reproducibility is a long-standing issue

- TOSCA description is used for logical preservation

- Docker images for tasks/workflows support physical preservation

- Changes tracking and automatic deployment also contribute to a comprehensive solution of the problem

- Integration of these techniques addresses majority of the issues related to workflow decay

# THANK YOU