

Building large-scale simulation pipelines using targets, Git and GitHub Actions

useR! 2024 Salzburg

Sergio Olmos

2024-07-11

About

- Statistician at Sanofi
- Member of [openstatsware](#)
- Innovative clinical trial designs and analytical methods
- Large scale simulations
- Software engineering best practices

sanofi

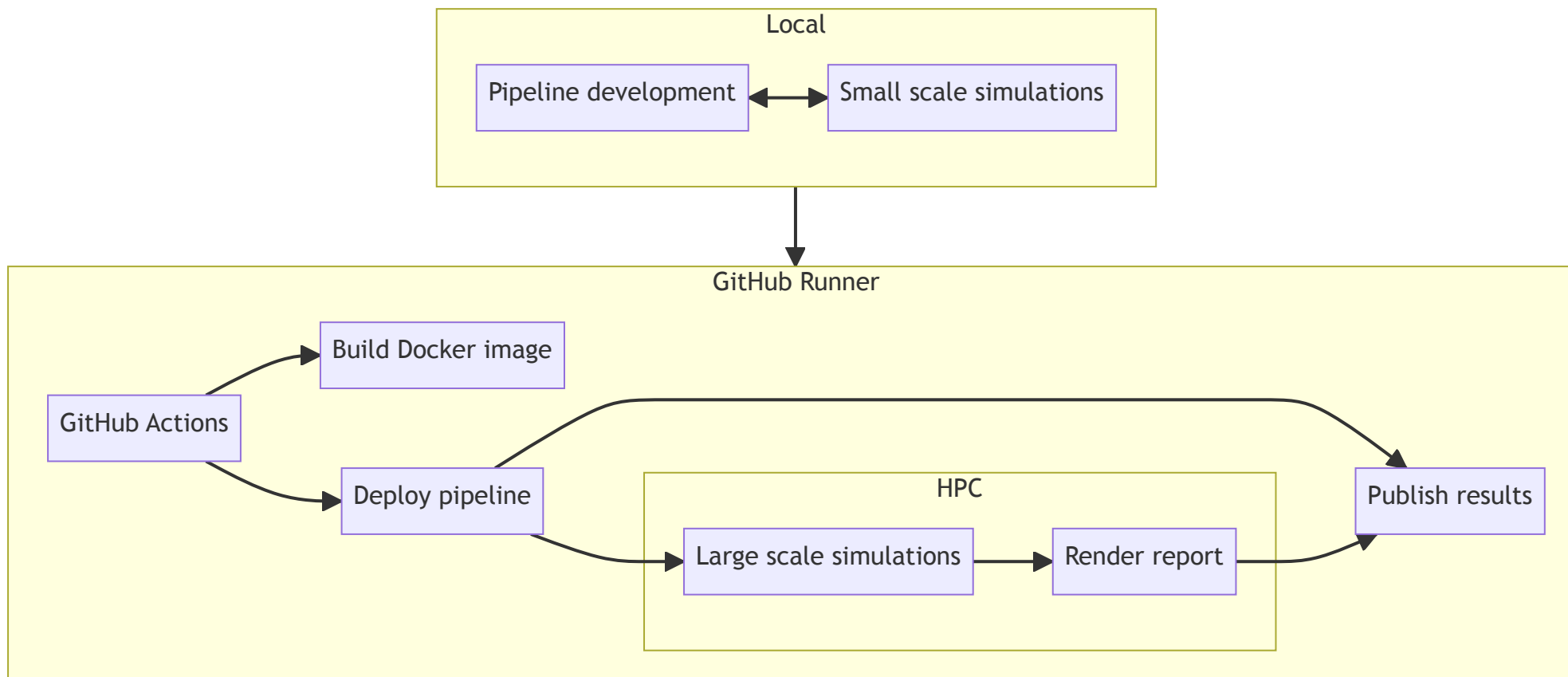


Overview

Tools

- {targets}
- Docker
- Git
- GitHub Actions

Pipeline



{targets}

Pipeline management

R-based

```
1 # _targets.R
2 library(targets)
3 tar_option_set(
4   packages = c("dplyr", "ggplot2")
5 )
6
7 list(
8   tar_target(
9     name = data_sim,
10    command = sim_data(reps = 100)
11  ),
12  tar_target(
13    name = model_fit,
14    command = fit_model(data_sim)
15  )
16 )
```

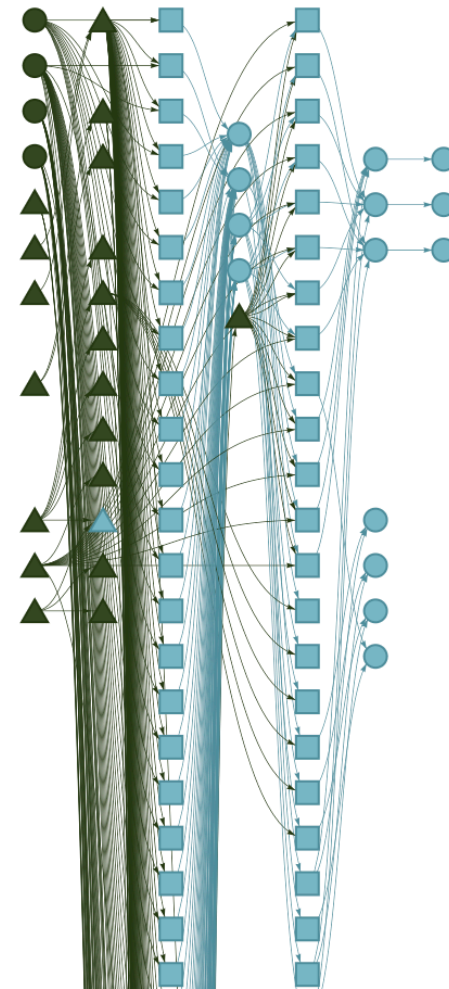
Run entire pipeline with:

```
1 targets::tar_make()
```


Parallelization via {crew}

```
1 # _targets.R
2 library(crew)
3 tar_option_set(controller = crew_controller_local(workers = 30))
```

- Dynamic branching
 - {tarchetypes}
- Distributed computing:
 - Local
 - SGE
 - SLURM
 - AWS Batch



Seamless testing-deployment workflow via {config}

config.yaml

```
1 default:
2   sims: 100
3   workers: 2
4 deployment:
5   sims: 10000
6   workers: 30
```

_targets.R

```
1 n_sims <- config::get("sims")
2 n_workers <- config::get("workers")
3 tar_option_set(crew_controller_local(workers = n_workers))
4 list(
5   tar_target(data_sim, sim_data(reps = n_sims))
6 )
```

Cloud computing

Cloud computing platform

- On-demand AWS EC2 instances
- Launch instance/cluster via API
- Custom Docker image
- Access to S3 buckets for persistent storage
- Choose HPC scheduler: SLURM, SGE, Kubernetes
- Live pipeline logs

Docker

Install all pipeline requirements

Production environment:

- Base images: [Rocker Project](#)
- System dependencies
- R version
- R packages

GitHub Actions

What are GitHub Actions?

- Workflow automation: CI/CD
- Runners:
 - GitHub-shared
 - Self-hosted
- Traceability
- Collaboration

Workflow

Before CI/CD

- Manually moving files
- Working in inefficient mounted folders
- Running costly cloud instance for long periods of time (mostly idle)
- Interactive testing in the command line
- Non-traceability of results

After CI/CD

Local development:

- Free
- Offline
- Your development environment
 - Desktop IDE
 - Git client
- Interactive testing

Automated pipeline deployment:

- No manual file movement
- Traceable
- Efficient use of costly cloud computing

Build image

- Trigger workflow only when changes in:
 - `Dockerfile`
 - Dependencies (e.g., `renv.lock`)
- Build image in GitHub runners
- Push image to cloud computing registry

Deploy pipeline

- Custom trigger:
 - `git push <deploy-branch>`
 - `git tag`
- Custom GitHub Action:
 - Specify Docker image
 - Specify HPC specs
 - Automated `git checkout`
 - Pipeline execution command

Deployment example

```
1 name: Deploy pipeline
2 on:
3   push:
4     branches: deploy
5 jobs:
6   deploy-pipeline:
7     runs-on: self-hosted
8     steps:
9       - name: Run pipeline
10         uses: org/custom-github-action/run-pipeline@main
11         with:
12           token: ${ secrets.API_TOKEN }
13           docker-image: user/custom-docker-image:latest
14           instance-type: r6i.32xlarge
15           command: sh _deploy.sh # targets::tar_make()
```

Publish results

After deployment workflow:

- targets pipeline includes Quarto/RMarkdown documents reporting results

```
1 tarchetypes::tar_quarto(report, "report.qmd")
```

- Publish rendered HTML to GitHub Pages/Posit Connect
- See results
- Share with collaborators

Questions?

✉ sergio.olmos.pardo@gmail.com

🐙 [solmos](#)

in [sergio-olmos](#)

✉ @solmos@fosstodon.org