

# FORCE on Nextflow: Scalable Analysis of Earth Observation data on Commodity Clusters

Fabian Lehmann<sup>1</sup>, David Frantz<sup>2,4</sup>,  
Sören Becker<sup>3</sup>, Ulf Leser<sup>1</sup>, Patrick Hostert<sup>4</sup>

<sup>1</sup> Institute for Computer Science, Humboldt-Universität zu Berlin

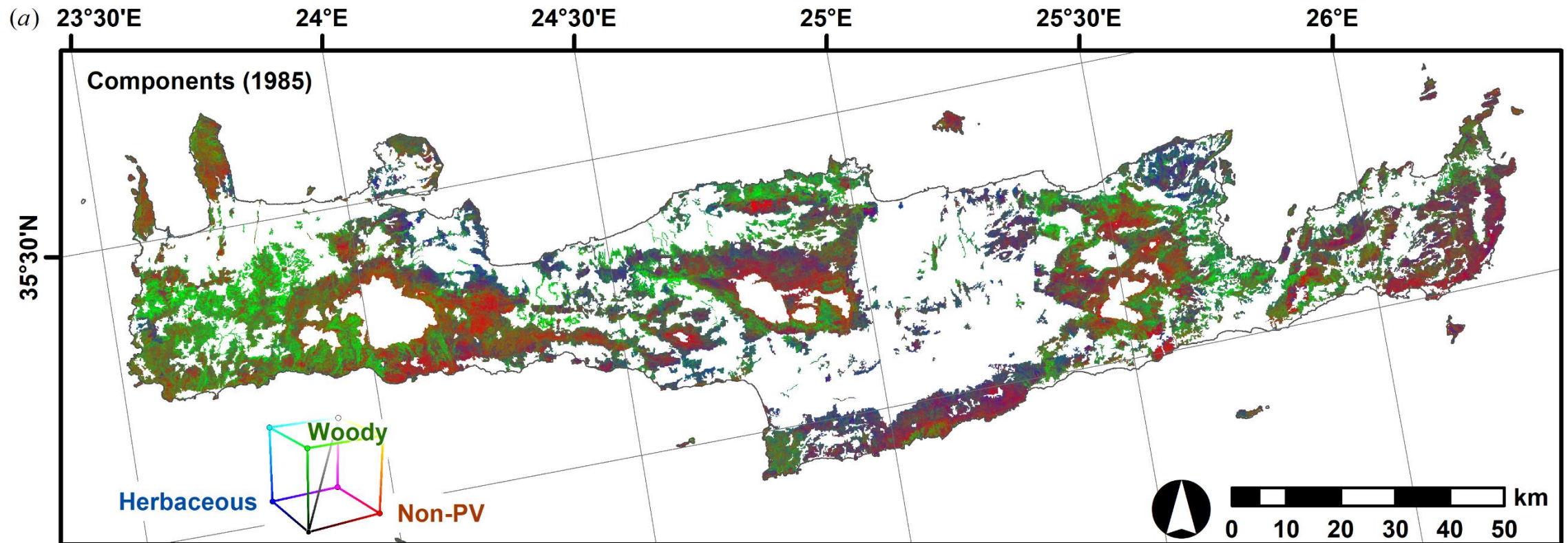
<sup>2</sup> Geoinformatics - Spatial Data Science, Universität Trier

<sup>3</sup> Distributed and Operating Systems, Technische Universität Berlin

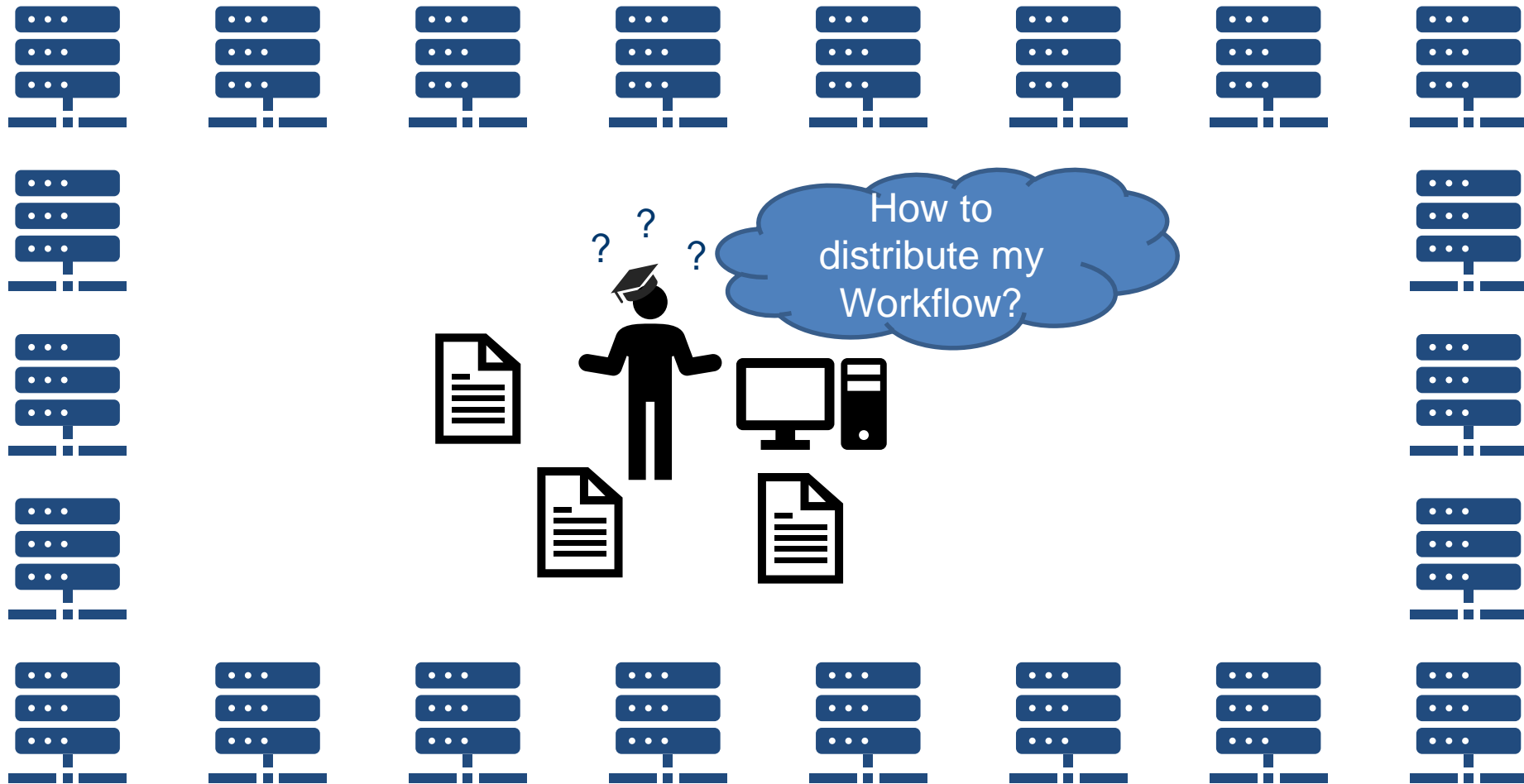
<sup>4</sup> Geography Department, Humboldt-Universität zu Berlin



## Crete: rangeland degradation analysis



# How to execute a distributed workflow



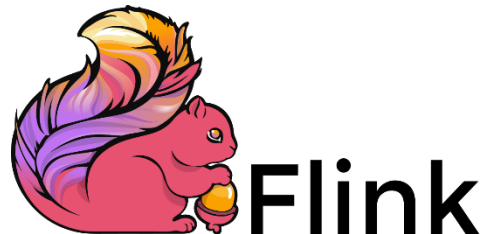
## Workflowsystems



kubernetes



## Workflowsystems



kubernetes



## Workflowsystems



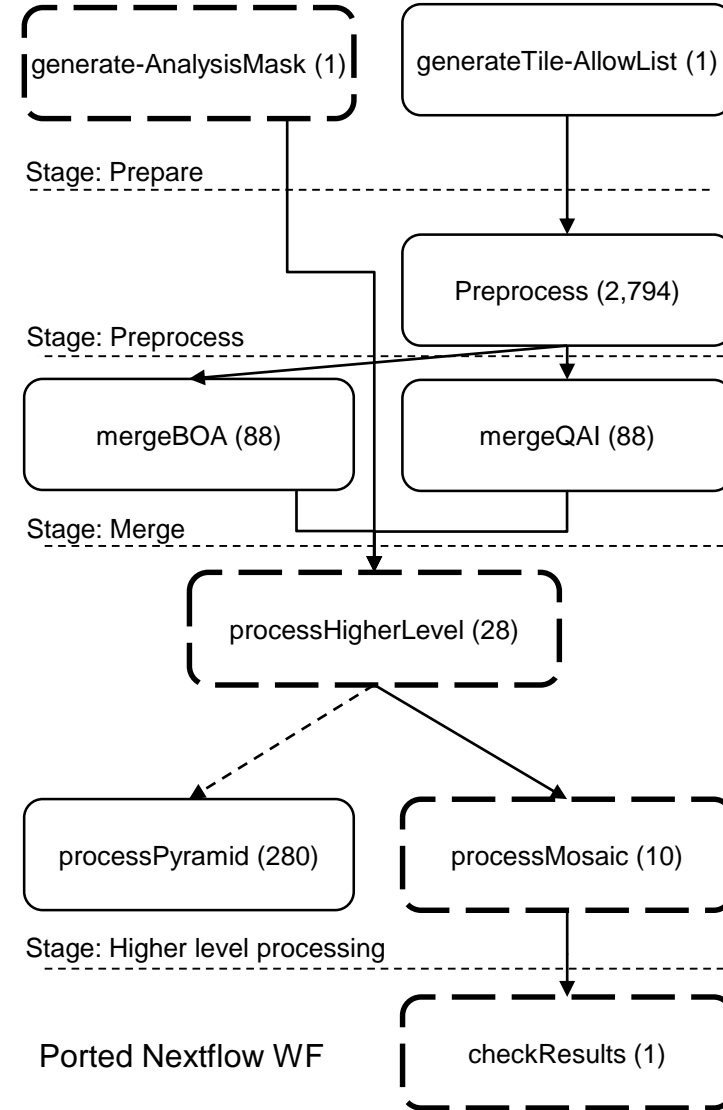
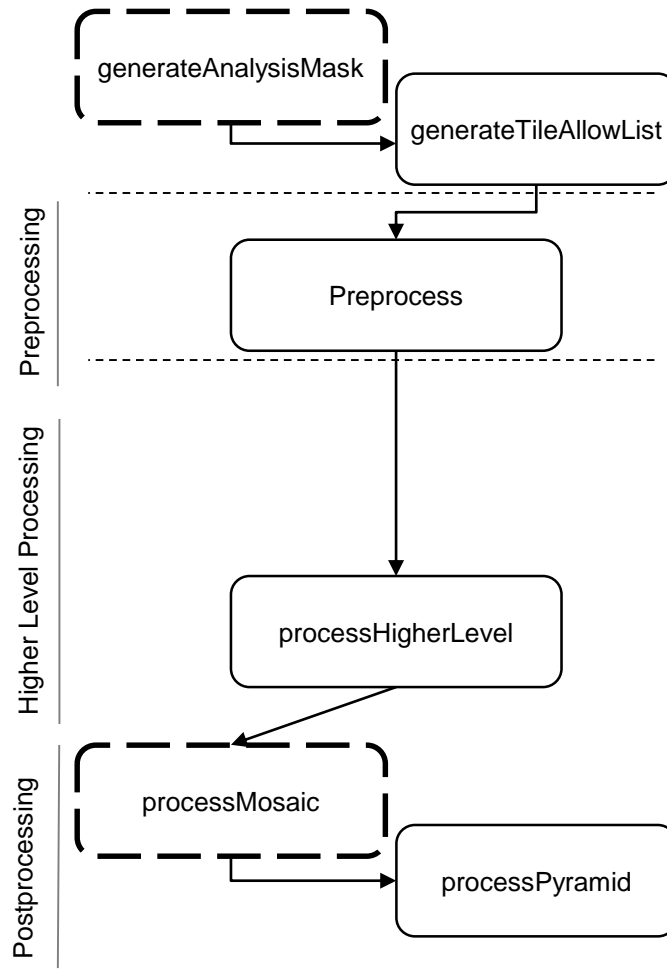
nextflow



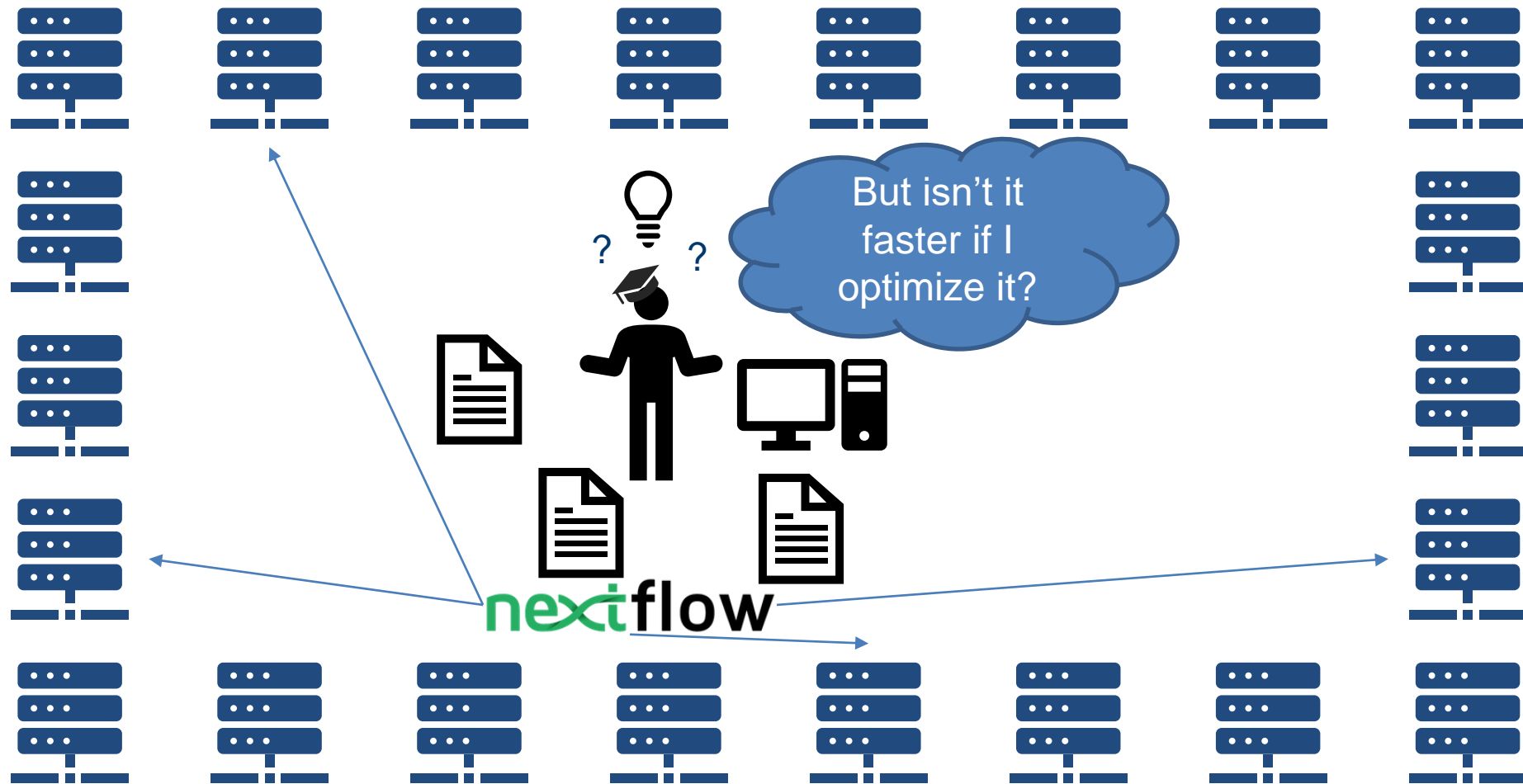
kubernetes



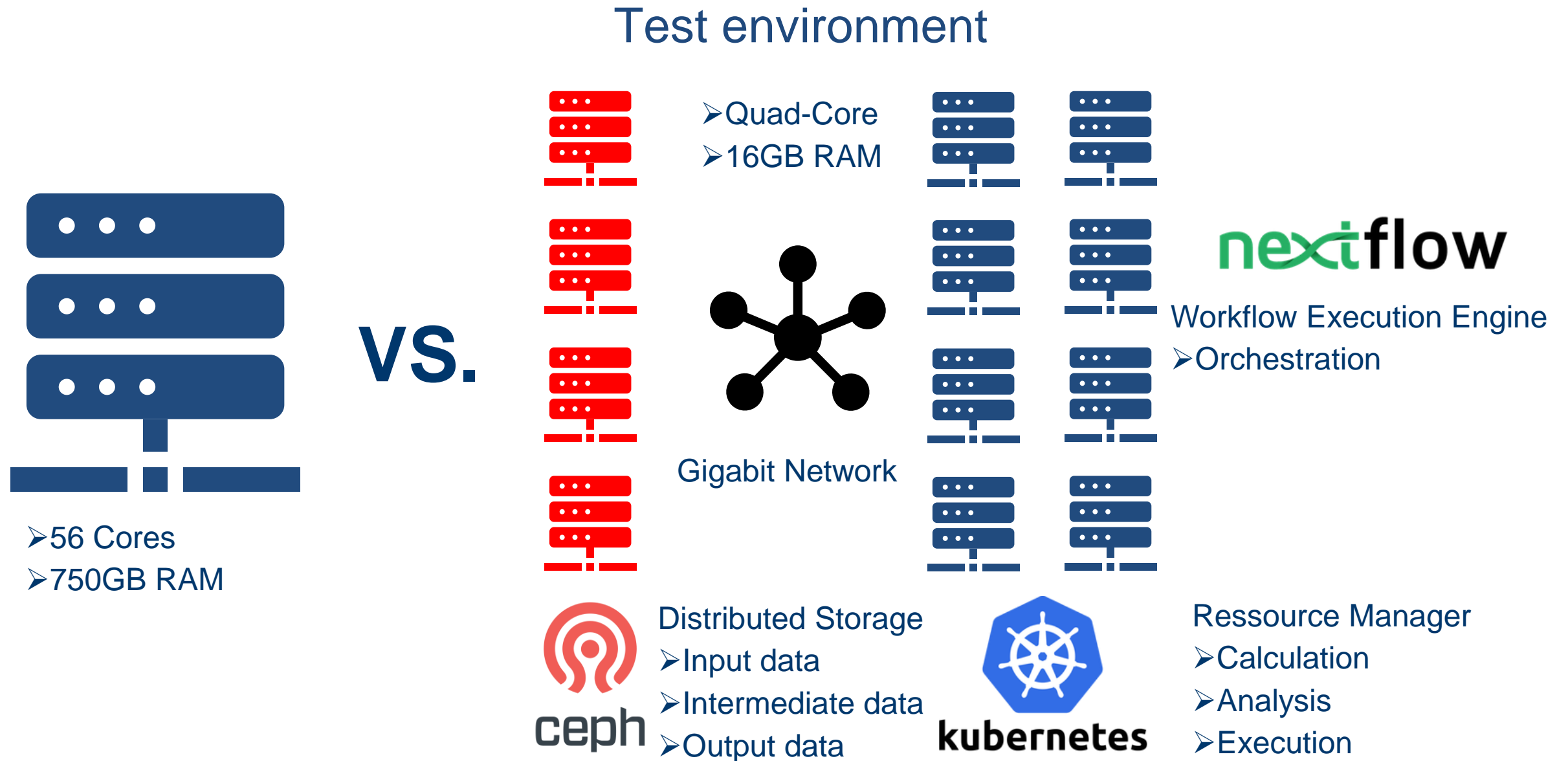
# How to adopt a BASH-Script to Nextflow



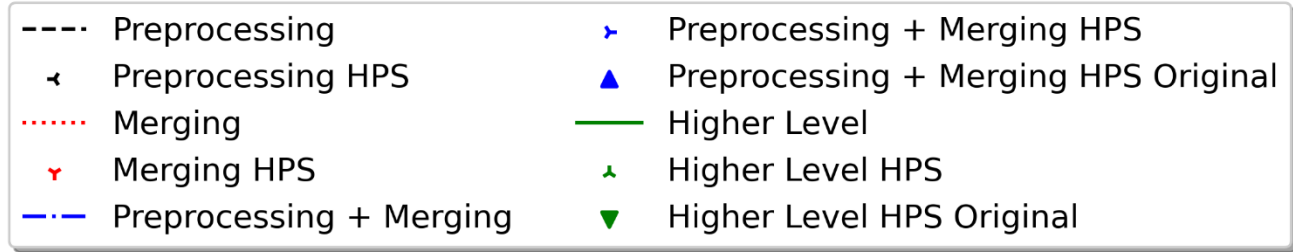
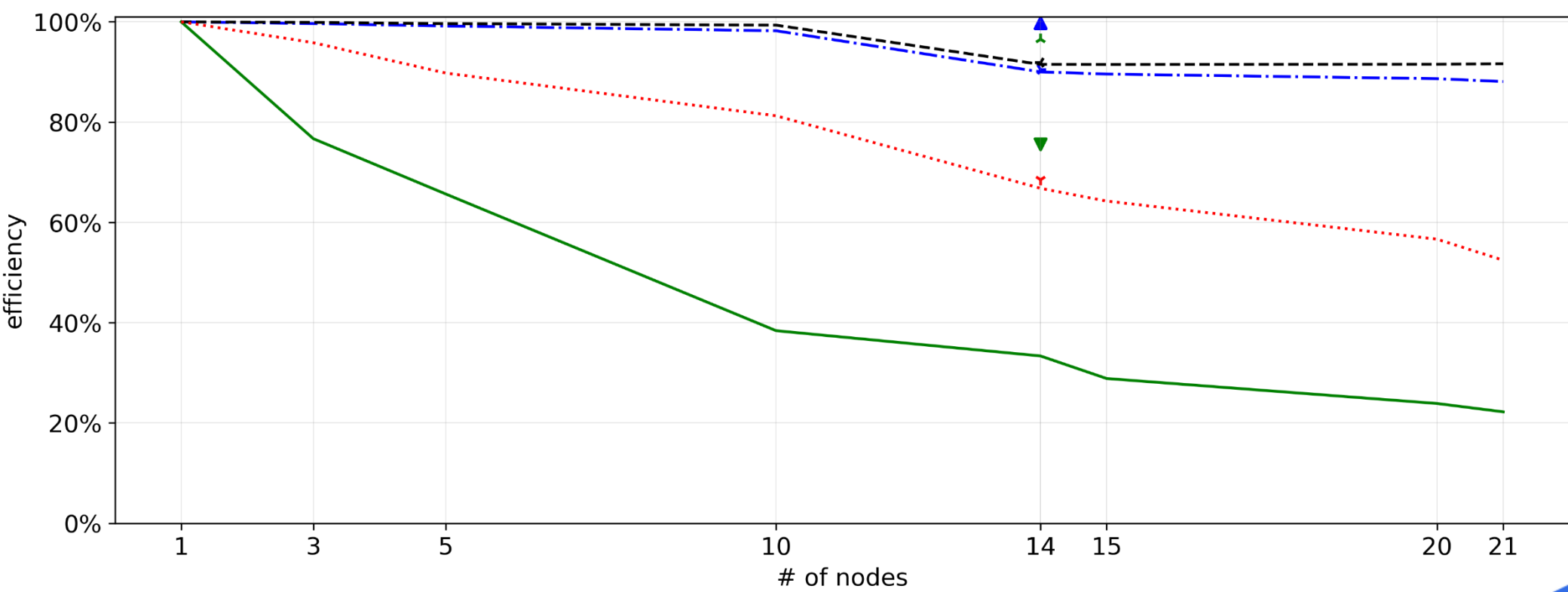
## How to run a distributed workflow







# Scalability analysis

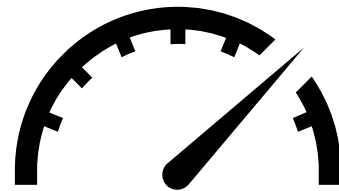
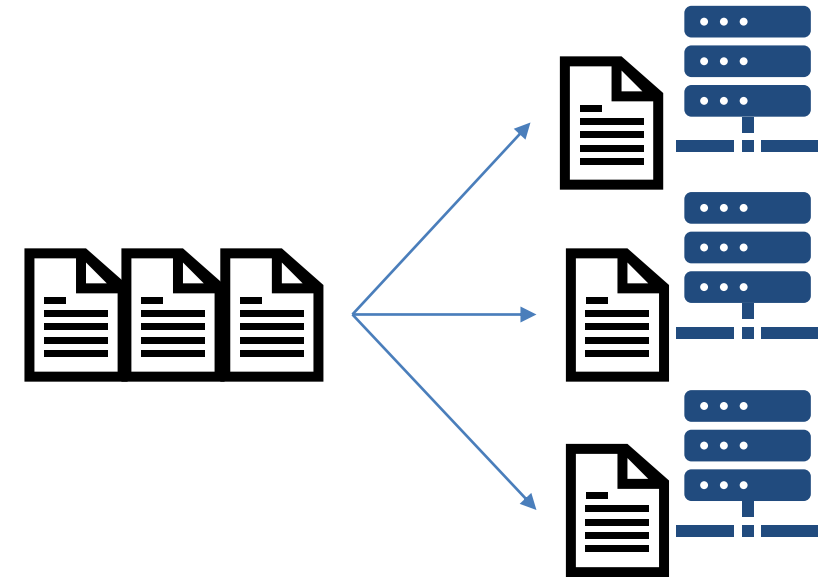


## Conclusion



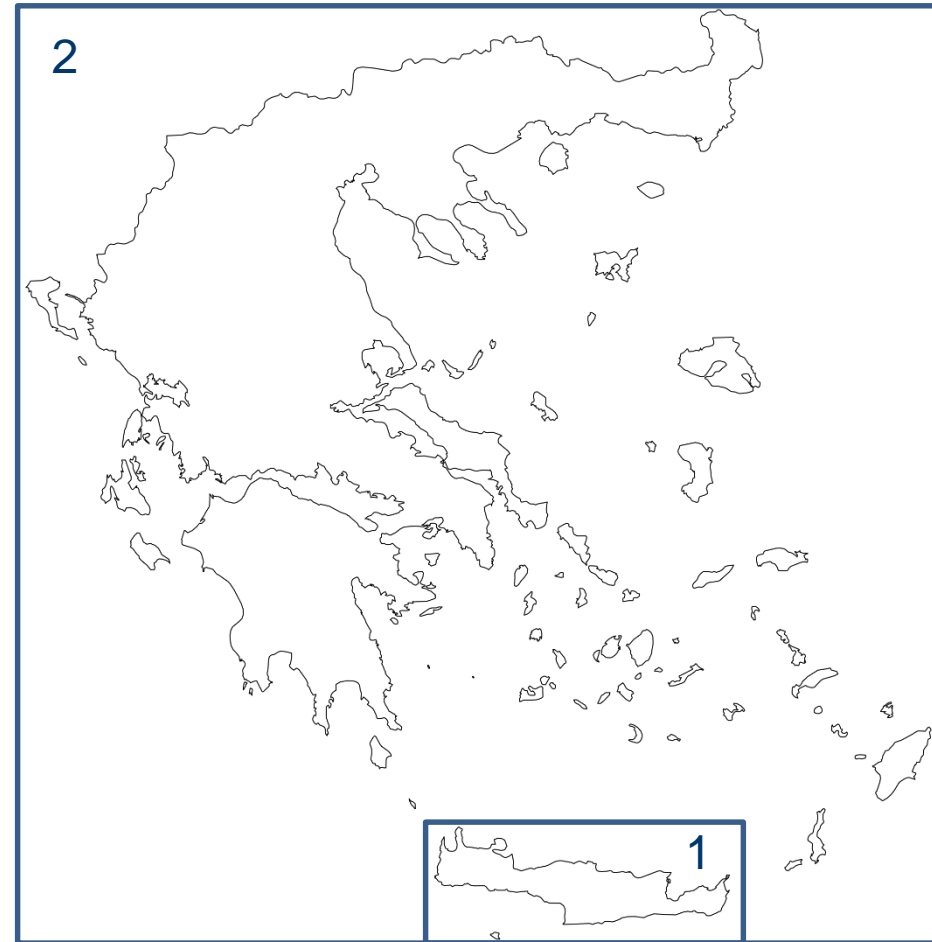
nextflow

- Learn additional workflow language
- Workflow Engine supports scalability
- I/O heavy tasks hardly scale
- Compute heavy tasks scale quite well
- Parallelism is defined implicitly



## Outlook

- Test scalability on larger cluster
- Run workflow for larger region
- Location aware scheduling



© Vemaps.com

**Contact:**

[fabian.lehmann@informatik.hu-berlin.de](mailto:fabian.lehmann@informatik.hu-berlin.de)

[david.frantz@uni-trier.de](mailto:david.frantz@uni-trier.de)

[soeren.becker@tu-berlin.de](mailto:soeren.becker@tu-berlin.de)

[leser@informatik.hu-berlin.de](mailto:leser@informatik.hu-berlin.de)

[patrick.hostert@geo.hu-berlin.de](mailto:patrick.hostert@geo.hu-berlin.de)

**Thank you  
for your attention!**

**FONDA – Foundations of Workflows for Large-Scale  
Scientific Data Analysis**

DFG Collaborative Research Center 1404 at Humboldt-Universität zu Berlin

**DFG** Deutsche  
Forschungsgemeinschaft



**FONDA**