

# FORCE on Argo Workflow

## Studienprojekt Exposé

Robin Knapp

4. März 2022

### 1 Einleitung

Workflows zur wissenschaftlichen Datenanalyse haben sich in unterschiedlichen wissenschaftlichen Disziplinen wie der Genom-[10] und der Klimaforschung [3] als Werkzeug etabliert. Sie verfügen über eine Systematik, welche die Ausführung mehrerer Einzelkomponenten, etwa einer umfangreichen Datenanalyse, verwaltet und koordiniert. Workflow stellen eine Umgebung dar, welche die Ausführung unabhängiger Computerprogramme steuert. Diese Computerprogramme stehen teilweise in Abhängigkeit zu einander, indem die Ausgabe der einen Anwendung durch die nächste weiterverarbeitet wird. Durch diese Abhängigkeiten formen sie einen gerichteten Graphen, dessen Knoten die einzelnen Computerprogramme bilden. Die Kanten des Graphen stellen die Abhängigkeiten durch Datenflüsse dar. Diese werden beispielsweise als Dateien zwischen den Knoten ausgetauscht [1]. Wissenschaftliche Workflows werden typischerweise von Wissenschaftlern für eine spezifische Systemlandschaft entwickelt. Außerhalb dieser Systemlandschaft ist der Workflow damit nicht ohne weiteres ausführbar, was die Reproduzierbarkeit von wissenschaftlichen Workflows einschränkt. Damit lassen sich auch die wissenschaftlichen Erkenntnisse nicht unabhängig von der verwendeten Systemlandschaft überprüfen bzw. reproduzieren[4].

Workflow Management Systeme implementieren umfangreiche Funktionalitäten, um Workflows auf Computer Ressourcen unabhängig von der gegebenen Systemlandschaft auszuführen. Hierzu gehören Strategien für die parallele Ausführung von Computerprogrammen, die Verteilung dieser auf gegebenen Ressourcen, das sogenannte Scheduling sowie die Verwaltung von Daten in verteilten Systemen [1]. Konkrete Workflow Management Systeme sind etwa Nextflow [7], Airflow [2] oder Pegasus [5]. Der Einsatz dieser Workflow Management Systeme erfordert jedoch eine Einarbeitung in das jeweilige System, um dessen Vorteile - etwa in Bezug auf die parallele Ausführung wissenschaftlicher Anwendungen - ausschöpfen zu können. Dieser Aufwand bedingt potentiell, dass bei der Entwicklung eines Workflows vom Einsatz eines Workflow Management Systems abgesehen wird [9]. Im Rahmen des Studienprojekts soll ebendiese Abwägung zwischen dem Aufwand und dem Nutzen eines Workflow Management Systems untersucht werden.

## 2 Zielsetzung

Ziel des avisierten Studienprojekts ist die Portierung eines Workflows auf ein Workflow Management System, welches es potentiell ermöglicht den wissenschaftlichen Workflow in heterogenen Umgebungen zu reproduzieren. Durch die Portierung soll zur Forschung im Bereich der Reproduzierbarkeit von wissenschaftlichen Workflows beigetragen werden. Hierbei soll speziell untersucht werden, welcher Aufwand bei der Portierung eines bestehenden Workflows auf ein Workflow Management System entsteht.

Vor dem Hintergrund, dass der betrachtete Workflow bereits durch Lehmann et al. [9] auf Nextflow portiert wurde, kann die Portierung auf ein zusätzliches Workflow Management System einen Beitrag zu weiterer Forschung im Bereich der wissenschaftlichen Workflows leisten. So kann die Portierung eines Workflows auf unterschiedliche Workflow Management Systeme die Grundlage für Untersuchungen im Bereich des Performance-Modeling [11] bilden und somit zu weiterer Forschung beitragen.

## 3 Forschungsstand

Lehmann et al. [9] untersuchen den Kompromiss zwischen Aufwand und Vorteilen bei der Portierung eines bestehenden Workflows auf ein Workflow Management System. Im Rahmen ihrer Untersuchung wurde ein auf FORCE basierender Workflow auf das Workflow Management System Nextflow portiert [9]. FORCE ist eine wissenschaftliche Anwendung für die Verarbeitung von Daten, welche anhand von Erdüberwachungssatelliten gesammelt werden. Hierzu kombiniert FORCE verschiedene Werkzeuge, die sich über entsprechende Konfigurationen zu einem Workflow verbinden lassen<sup>1</sup>. Diese Werkzeuge laden Daten für die gewählten Bildausschnitte aus entsprechenden Datenbanken herunter und bereiten diese in mehreren Schritten für eine weitere Auswertung vor [8]. Lehmann et al. [9] stellen fest, dass die Entwicklung eines Workflows auf Nextflow potentiell für umfangreichere Analyseaufgaben sinnvoll sein kann. Sie stellen jedoch auch fest, dass speziell in Bezug auf die Ausführung in verteilten Umgebungen, Inkompatibilitäten zwischen den FORCE Werkzeugen mit den Konzepten von Nextflow bestehen. Dieses führt in der Folge dazu, dass teilweise die FORCE spezifische Logik übernommen werden musste. Die Autoren stellen zusätzlich in Aussicht, die Untersuchung mit dem Workflow Management System Apache Airflow zu wiederholen [9]. Vor diesem Hintergrund drängt sich - wie im vorliegenden Studienprojekt avisiert - eine Portierung desselben Workflows auf ein weiteres Workflow Management System auf, um so einen umfassenden Vergleich zwischen den verschiedenen Ansätzen zu ermöglichen.

Dessalk et al. [6] untersuchen einen Ansatz zur Ausführung von Big Data Workflows. Im Rahmen ihrer Untersuchung vergleichen sie relevante Workflow

---

<sup>1</sup> <https://github.com/CRC-FONDA/FORCE2NXF-Rangeland#original-workflow>

Management Systeme. Von diesen Workflow Systemen nutzen sie Argo Workflows (Argo) <sup>2</sup>, welches auf Kubernetes aufbaut, um dieses mit einer von ihnen entwickelten Workflow zu vergleichen. Die Autoren wählen Argo, da es über keine integrierte Unterstützung für Nebenläufigkeit verfügt, stellen das Workflow Management System unabhängig davon aber als relevantes System dar [6]. In der Folge bietet sich Argo als relevante Option für die Portierung von FORCE an.

## 4 Vorgehen

Der bereits von Lehmann et al. [9] untersuchte Workflow zur Satellitenbild-Analyse wird von einem proprietären System auf ein weiteres Workflow Management System portiert. Als Workflow Management System wird Argo genutzt, welches die definierten Einzelschritte des Workflows auf Kubernetes ausführt. Durch die Ausführung auf einem Kubernetes Cluster ist grundsätzlich eine Vergleichbarkeit zwischen den beiden Ansätzen gegeben bzw. können diese potentiell auf dem selben Cluster ausgeführt werden. Hierbei sollen zunächst die in FORCE verwendeten Werkzeuge einzeln in Containern ausgeführt werden. Anschließend wird die Ausführung der Komponenten mit Argo umgesetzt. Die Definition der einzelnen Schritte des Workflows erfolgt in Argo als sogenannte Tasks. Dieses erfolgt über eine Domain Specific Language (DSL) in YAML Dateien. Jene Einzelschritte werden dann als Graph in Argo definiert, sodass der Workflow in der angestrebten Reihenfolge abläuft. Zudem muss die Bereitstellung von Daten zwischen den Einzelschritten des Workflows abgebildet werden. Hierfür sieht Argo das Konzept von sogenannten Artifacts vor.

Im Anschluss an die Modellierung des Workflows in der Argo spezifischen DSL wird der Workflow auf einem Kubernetes Cluster ausgeführt. Um eine Vergleichbarkeit der Portierung zu den Experimenten von Lehmann et al. [9] zu erreichen, wird der Workflow mit den selben Daten ausgeführt. Verwendet wird hierbei ein Bildausschnitt der griechischen Insel Kreta <sup>3</sup>.

Im Rahmen des Experiments wird der Workflow auf einem Kubernetes Cluster ausgeführt und die benötigte Zeit für die vollständige Abarbeitung des Workflows gemessen. Dieses Experiment wird mit einer unterschiedlich Zahl Cluster-Knoten durchlaufen. Ziel ist es die Zeit zwischen dem Beginn der Ausführung bis zur vollständigen Abarbeitung des Workflows im Verhältnis zur Anzahl der verwendeten Cluster-Knoten zu erfassen.

## Literatur

- [1] R. M. Badia Sala, E. Ayguadé Parra, and J. J. Labarta Mancho. Workflows for science: a challenge when facing the convergence of hpc and big data. *Supercomputing frontiers and innovations*, 4(1):27–47, 2017.

<sup>2</sup> <https://argoproj.github.io/argo-workflows/>

<sup>3</sup> <https://github.com/CRC-FONDA/FORCE2NXF-Rangeland#area-of-interest-vector>

- [2] M. Beauchemin. Airflow: a workflow management platform. <https://medium.com/airbnb-engineering/airflow-a-workflow-management-platform-46318b977fd8>, 2015.
- [3] J. Buchner, H. Yin, D. Frantz, T. Kuemmerle, E. Askerov, T. Bakuradze, B. Bleyhl, N. Elizbarashvili, A. Komarova, K. E. Lewińska, et al. Land-cover change in the caucasus mountains since 1987 based on the topographic correction of multi-temporal landsat composites. *Remote Sensing of Environment*, 248:111967, 2020.
- [4] S. Cohen-Boulakia, K. Belhajjame, O. Collin, J. Chopard, C. Froidevaux, A. Gaignard, K. Hinsén, P. Larmande, Y. Le Bras, F. Lemoine, et al. Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities. *Future Generation Computer Systems*, 75: 284–298, 2017.
- [5] E. Deelman, J. Blythe, Y. Gil, C. Kesselman, G. Mehta, S. Patil, M.-H. Su, K. Vahi, and M. Livny. Pegasus: Mapping scientific workflows onto the grid. In *European Across Grids Conference*, pages 11–20. Springer, 2004.
- [6] Y. D. Dessalk, N. Nikolov, M. Matskin, A. Soyly, and D. Roman. Scalable execution of big data workflows using software containers. In *Proceedings of the 12th International Conference on Management of Digital EcoSystems*, pages 76–83, 2020.
- [7] P. Di Tommaso, M. Chatzou, E. W. Floden, P. P. Barja, E. Palumbo, and C. Notredame. Nextflow enables reproducible computational workflows. *Nature biotechnology*, 35(4):316–319, 2017.
- [8] D. Frantz. Force—landsat+ sentinel-2 analysis ready data and beyond. *Remote Sensing*, 11(9):1124, 2019.
- [9] F. Lehmann, D. Frantz, S. Becker, U. Leser, and P. Hostert. Force on nextflow: Scalable analysis of earth observation data on commodity clusters. In *Int. Workshop on Complex Data Challenges in Earth Observation*, 2021.
- [10] C. Schiefer, M. Bux, J. Brandt, C. Messerschmidt, K. Reinert, D. Beule, and U. Leser. Portability of scientific workflows in ngs data analysis: a case study. *arXiv preprint arXiv:2006.03104*, 2020.
- [11] C. Witt, M. Bux, W. Gusew, and U. Leser. Predictive performance modeling for distributed batch processing using black box monitoring and machine learning. *Information Systems*, 82:33–52, 2019.