

Workflow and Workload Management
with Airflow & Ray

Stéphane Bisinger

29th January 2026, Zürich



Stéphane Bisinger

- Software Engineer
- Worked at Homegate, siroop, and many more
- Web development: Backend, Frontend, Fullstack
- Infrastructure: Serverless Architectures, Kubernetes
- Extreme Programming, TDD, Continuous Delivery
- At KOF since 2022
 - Data import
 - Coding support and architecture
 - Infrastructure development and maintenance

KOF Swiss Economic Institute

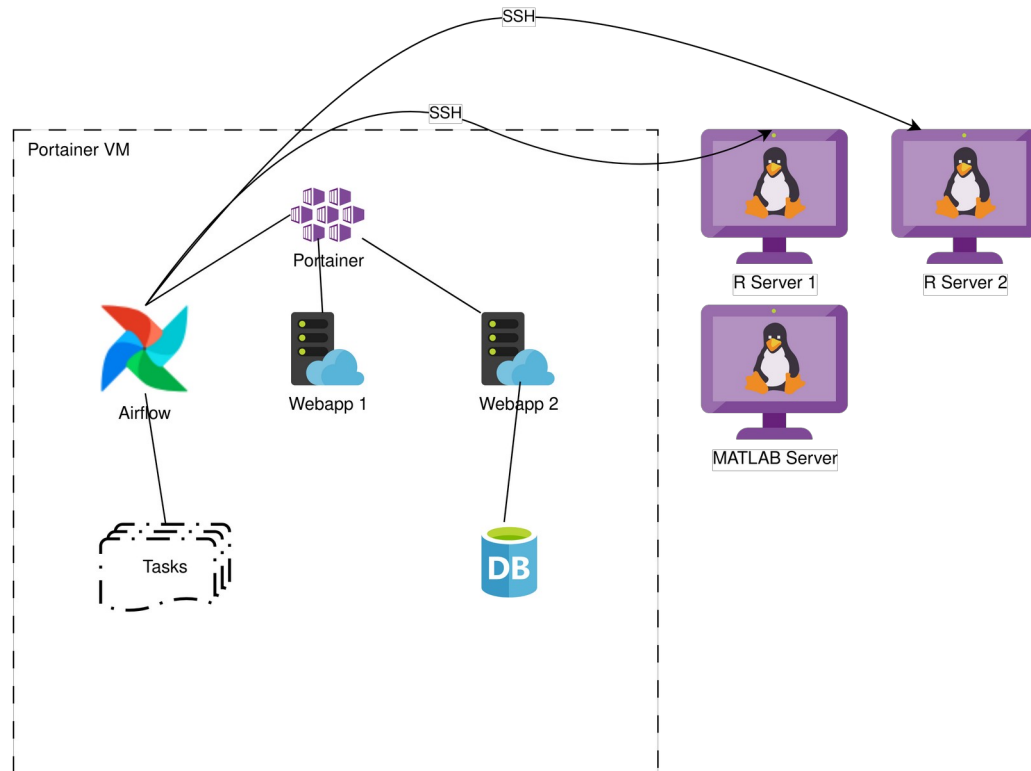
- The KOF sees itself as a **bridge builder** between **research and society** and as an **exchange platform for economists**
- Fulfills so-called **national tasks** that have been assigned to it by the **federal government**
- Some of the main products from KOF are:
 - Quarterly updated macroeconomic forecasts for the Swiss economy. (R)
 - Business surveys of Swiss companies. (R)
 - The Nowcasting Lab <https://nowcastinglab.org> (MATLAB, Python)
 - ...and many other Surveys, Forecasts, and Indicators <https://kof.ethz.ch>

1. VM-base infrastructure
2. Kubernetes
3. Ray

VM-based architecture



- Data import
 - Read an API or website to check if new data we are interested in is available
 - Download the data and parse the format
 - Save it in our internal database
- Seasonal Adjustment
 - When data is updated, calculate and store the seasonally adjusted variant of a time series
- Data aggregations
- Run prediction models
 - Load the required data
 - Calculate the prediction with Fancy Models
 - Store the output of the model



- There's a few VMs dedicated to running R and MATLAB programs
- A VM would host the docker host, with portainer as a UI
- Various web apps and other applications would run here as containers
- Airflow also runs as a container
 - It will launch simple R tasks as a standalone container
 - It will connect to other specific VMs for complex tasks requiring more computing resources

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
AI Agent Based Forecasting	VJA	221 / 30	0 12 * * *	2026-01-17, 11:00:00	2026-01-24, 11:00:00	4
Forschungsmonitoring sync authors	airflow	2	1 day, 0:00:00	2025-02-24, 15:31:31		1
Ray_Taskflow_Example_Existing_Cluster	airflow	5 / 9	None	2025-11-21, 15:26:45		3 / 2
bts_agg_k8s	KOF	534 / 22	None	2026-01-24, 15:00:06		31
bts_agg_k8s_test	KOF	87 / 17	None	2026-01-27, 01:14:38		3
bts_agg_test_trigger	KOF	3	None	2026-01-27, 01:13:55		7
bts_agg_test_trigger_manual	KOF	20	None	2026-01-20, 17:22:47		1
bts_agg_total_k8s	KOF	87 / 11	None	2026-01-24, 19:36:52		11
bts_agg_total_k8s_test	KOF	3 / 3	None	2025-10-08, 12:10:23		6
bts_agg_total_trigger_manual_test	KOF	8	None	2023-08-21, 19:01:21		1
bts_agg_trigger_pending_k8s	KOF	534 / 997	* * * * *	2026-01-28, 17:52:00	2026-01-28, 13:32:00	2
bts_agg_trigger_pending_k8s_test	KOF	890 / 46	* * * * *	2025-09-24, 15:31:00	2026-01-28, 13:31:00	2
bts_customer_email	KOF	38 / 20	None	2025-10-24, 13:11:32		1

DAG: swissdata_modular Download and process data from Swiss public institutions

Schedule: 0,31 8,9,17 * * * * * Next Run ID: 2026-01-28, 16:31:00 UTC

28 / 01 / 20 17:53:02 All Run Types All Run States Clear Filters

Press shift + / for Shortcuts

Duration: 00:01:29 to 00:00:00

Tasks: latest_only, find_updated_swissdata_sets, trigger_updated_sets, ch-fso_to_csv, ch-adecco_to_csv, ch-seco_to_csv, ch-snb_to_csv, ch.snb.amarbma-to_csv, ch.snb.amarbma-source.staging, ch.snb.amarbma-source.prod, ch.snb.ambeschkla-to_csv, ch.snb.ambeschkla-source.staging, ch.snb.ambeschkla-source.prod, ch.snb.ausschwarm-to_csv, ch.snb.ausschwarm-source.staging, ch.snb.ausschwarm-source.prod, ch.snb.auvekomq-to_csv, ch.snb.auvekomq-source.staging, ch.snb.auvekomq-source.prod, ch.snb.auvercurrq-to_csv, ch.snb.auvercurrq-source.staging

Airflow is cron on steroids

Run tasks on a schedule. Keep track of runs and their status and logs.

- Interval-based scheduling
- Event-based scheduling
- Custom timetables
- ...and more!

```
Dockerfile
1 FROM rocker/r-ver:3.6.0 as deps
2
3 RUN apt-get update && apt-get install -y \
4     libpq-dev \
5     libcurl4-openssl-dev \
6     libxml2-dev \
7     ...
8
9 RUN mkdir /tasks
10 COPY tasks/ /tasks/
11
12 # Installed internal packages used by Airflow scripts
13 RUN install2.r --repos registry.kof.ethz.ch kofpkg1 kofpkg2 kofpkg3
14
15
16 FROM deps
17
18 COPY .Rprofile /home/kofbaser/.Rprofile
19 COPY entrypoint.sh /entrypoint.sh
20 RUN chmod +x /entrypoint.sh
21
22 WORKDIR /home/kofbaser
23
24 ENTRYPOINT [ "/entrypoint.sh" ]
```

```
entrypoint.sh
1 #!/bin/bash
2
3 if [[ ! -f /tasks/$1.R ]]; then
4     echo not a task file, doing whatever
5     exec $@
6 else
7     echo running /tasks/$1.R ${@: 2}
8     Rscript /tasks/$1.R ${@: 2}
9 fi
```

```
swissdata.py

1 from datetime import datetime
2 from airflow import DAG
3 from airflow.models import Variable
4 from airflow.operators.docker_operator import DockerOperator
5
6 default_args = {
7     'start_date': datetime(2017, 3, 20)
8 }
9
10 dag = DAG('swissdata', description='Run swissdata',
11         schedule_interval = '0,31 8,9,17 * * *',
12         default_args = default_args, catchup=False)
13
14 pg_password = Variable.get("docker_pg_password")
15
16 with dag:
17     source_swissdata = DockerOperator(
18         task_id = 'source_swissdata',
19         image = 'registry.ethz.ch/kof/kofdocker/kofbase:0.2.0',
20         user='kofbaser',
21         volumes = [
22             'swissdata-output:/swissdata',
23             'kofbase-tasks:/tasks'
24         ],
25         docker_url = 'unix:///var/run/docker.sock',
26         environment = {
27             'PG_PASSWORD': pg_password,
28             'PG_HOST': 'db.kof.ethz.ch'
29         },
30         command = 'source_swissdata'
31     )
32
```

```
business_tendency_survey.py

1 import datetime
2 from airflow import DAG
3 from airflow.contrib.operators.ssh_operator import SSHOperator
4
5
6
7 default_args = {
8     'start_date': datetime.datetime(2019, 11, 12)
9 }
10
11 dag = DAG(
12     'bts_agg_prod',
13     description='Runs business tendency survey aggregation [PROD]',
14     schedule_interval = None,
15     default_args = default_args,
16     catchup = False
17 )
18
19 with dag:
20     aggregate = SSHOperator(
21         task_id='aggregate',
22         ssh_conn_id='rcalc_prod',
23         command="""Rscript -e 'library(kofbts); runAggregation( \
24             "{{dag_run.conf.survey_name}}", \
25             {{dag_run.conf.run_of_analysis}}, \
26             {{dag_run.conf.year}}, \
27             {{dag_run.conf.month}})'"""
28     )
29
```

Pros

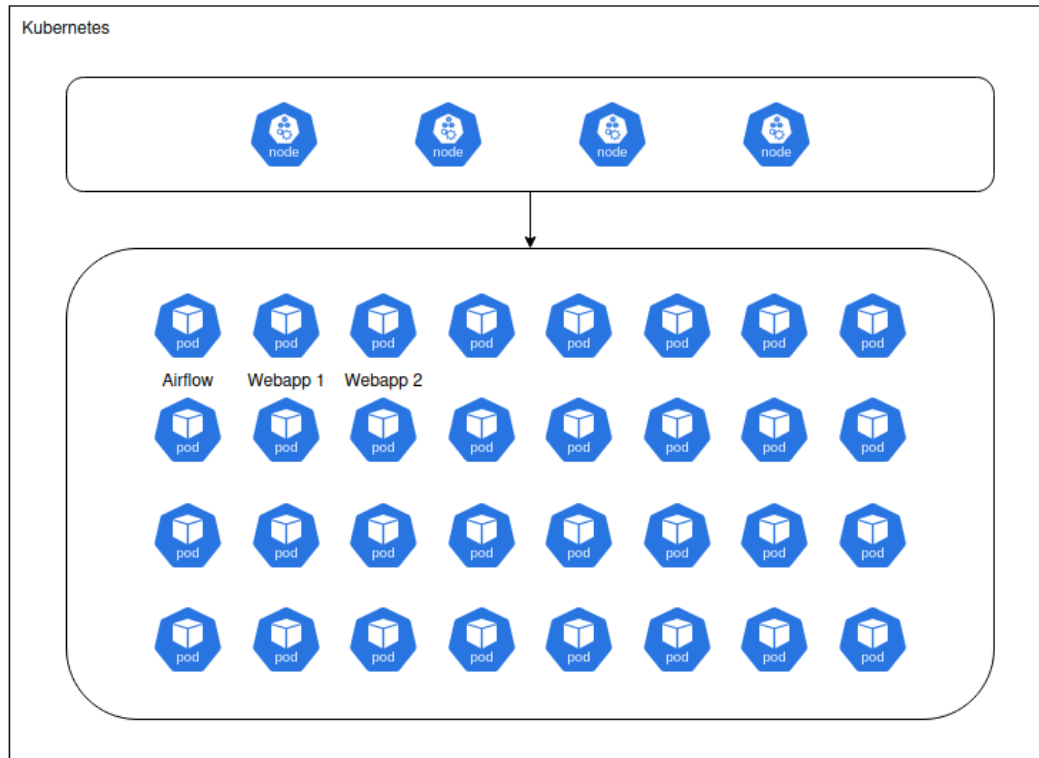
- Computations move away from personal computers, onto dedicated infrastructure that doesn't require a user sitting in front of the computer
- With Airflow, regularity and timing can be automated.
- Automatic notification in case of errors
- Log history of runs

Cons

- VMs are always running, whether they are needed or not
- Servers need to be installed with all necessary dependencies (and managed)
- Favors big scripts that do it all, rather than modular and parametrizable tasks
- Parallelism is not really encouraged
- Infrastructure management is fragmented and complex

Kubernetes





- The architecture is simply a kubernetes cluster: no special cases, Kubernetes manages workloads, routing, scaling, etc.
- The VMs are now just Kubernetes Nodes
- Everything is containerized: no specialized servers that require specific libraries or tooling.
- Airflow tasks are now run as pods

Running R tasks – now with renv!

```
Dockerfile
1 ARG R_VERSION=4.3.1
2 FROM docker.io/rocker/r-ver:${R_VERSION}
3
4 # libssl-dev must be installed first to work in older R versions. It might get
5 # downgraded to the correct version with the follow-up command.
6 RUN apt-get update \
7     && apt-get install -y libssl-dev \
8     && apt-get install -y \
9     libpq-dev \
10    libcurl4-openssl-dev \
11    libxml2-dev \
12    ... \
13    && rm -rf /var/lib/apt/lists/*
14
15 # Only renv needed
16 RUN install2.r renv --repos https://packagemanager.posit.co/cran/latest
17
18 # renv volume must be mounted to /renv in container for renv cache
19 ENV RENV_PATHS_PREFIX_AUTO=TRUE
20 ENV RENV_PATHS_CACHE=/renv/cache
21
22 COPY entrypoint.sh /entrypoint.sh
23 COPY run_task.R /
24 COPY Rprofile /home/kofbaser/.Rprofile
25 RUN chmod +x /entrypoint.sh
26
27 WORKDIR /home/kofbaser
28
29 ENTRYPOINT [ "/entrypoint.sh" ]
30
```

```
entrypoint.sh
1 #!/bin/bash
2 if [ "$1" = "run_task" ]; then
3     echo running task $3 in package $2
4     Rscript /run_task.R $2 $3
5 elif [[ ! -f /tasks/$1/run.R ]]; then
6     echo not a task file, doing whatever
7     exec $@
8 else
9     echo running /tasks/$1/run.R
10    cd /tasks/$1
11    R -e "renv::restore(); source('run.R')"
12 fi
```

```
run_task.R
1 args <- commandArgs(trailingOnly = TRUE)
2 package_name <- args[[1]]
3 command <- args[[2]]
4
5 o <- renv::refresh()
6 renv::install(package_name)
7 tasks_dir <- system.file("tasks", package = package_name)
8 # packages can have their own tasks in inst/tasks/<task>.R
9 script <- paste0(tasks_dir, "/", command, ".R")
10 message(sprintf("Executing script: %s", script))
11 sys.source(script, envir = topenv())
```

```
swissdata.py
1 set_name = "ch.fso.frv"
2 environment = "prod"
3 db_hosts = {"prod": "db.kof.ethz.ch"}
4 pg_password = Variable.get("prod_db_password")
5
6 with dag:
7     source_swissdata = KubernetesPodOperator(
8         namespace="airflow-worker",
9         image="registry.kof.ethz.ch/images/kofrenv:latest",
10        task_id=f"{set_name}-source.{environment}",
11        in_cluster=True,
12        get_logs=True,
13        env_vars=k8s_env_vars(
14            {
15                "PG_PASSWORD": pg_password,
16                "PG_HOST": db_host(environment),
17                "SWISSDATA_SET": set_name,
18            }
19        ),
20        volume_mounts=k8s_volume_mounts({"renv-cache": "/renv"}),
21        volumes=k8s_volumes({"renv-cache": "/renv"}),
22        arguments=["run_task", "kofpkg1", "source_swissdata"],
23    )
```

```
business_tendency_survey.py
1 pg_password = Variable.get("prod_db_password")
2
3 with dag:
4     aggregate = KubernetesPodOperator(
5         namespace="airflow-worker",
6         task_id="aggregate",
7         image="registry.kof.ethz.ch/images/kofrenv:4.5.0-1",
8         in_cluster=True,
9         get_logs=True,
10        env_vars=k8s_env_vars(
11            {
12                "pg_password": pg_password,
13                "nr_cores": "4",
14                "survey": "{{dag_run.conf.survey_name}}",
15                "run": "{{dag_run.conf.run_of_analysis}}",
16                "year": "{{dag_run.conf.year}}",
17                "month": "{{dag_run.conf.month}}",
18            }
19        ),
20        volume_mounts=k8s_volume_mounts({"renv-cache": "/renv"}),
21        volumes=k8s_volumes({"renv-cache": "/renv"}),
22        container_resources=k8s.V1ResourceRequirements(
23            requests={"memory": "24Gi", "cpu": "4000m"},
24            limits={"memory": "48Gi", "cpu": "12000m"},
25        ),
26        arguments=["run_task", "kof_business_tendency_survey", "aggregate"],
27    )
```

Pros

- Infrastructure is now standardized and managed through Helm Charts (Infrastructure-as-Code)
- Everything can be scaled automatically on-demand
- Monitoring everything is much simpler
- All the R-code is in the R package, not spread over multiple projects.
- Usage is simplified
- Encourages splitting into smaller modular tasks
- Can run any type of workload with the same architecture

Cons

- Kubernetes management is a complex skill that requires time to acquire
- Resource limits must be set to avoid greedy tasks taking over all the available resources
- Small tasks have a relatively high overhead

Ray



JOBS

Auto Refresh:

Request Status: Fetched jobs

Job List

Job ID Submission ID Page Size 10 Per Page Status

< 1 2 3 4 5 ... 54 >

Job ID	Submission ID	Entrypoint	Status	Status message	Duration	Tasks ?	Action
(no ray driver)	raysubmit_sAUeAN3NB3MVBKDK	Rscript /run_task.R kofsourcing source_alfred	SUCCEEDED	Job finished successfully. Expand	5m 21s	<div style="width: 100%;"></div>	Log
(no ray driver)	raysubmit_Vxuw8N13qxCkvvAr	Rscript /run_task.R kofsourcing source_macrobond	SUCCEEDED	Job finished successfully. Expand	3m 50s	<div style="width: 100%;"></div>	Log

Ray is an open source unified framework for scaling AI and Python applications.

It provides a simple, universal API for building distributed applications that can scale from a laptop to a cluster.

What's Ray?

Ray simplifies distributed computing by providing:

- **Scalable compute primitives:** Tasks and actors for painless parallel programming
- **Specialized AI libraries:** Tools for common ML workloads like data processing, model training, hyperparameter tuning, and model serving
- **Unified resource management:** Seamless scaling from laptop to cloud with automatic resource handling

```
ray.yaml
1 head:
2   rayVersion: {{ $vars.RAY_VERSION }}
3   enableInTreeAutoscaling: true
4   autoscalerOptions:
5     idleTimeoutSeconds: 300
6     imagePullPolicy: Always
7 worker:
8   image:
9     repository: registry.kof.ethz.ch/images/ray-kofrenv
10    tag: {{ $tag }}
11  rayStartParams:
12    resources: '{"r": 1}'
13  replicas: 0
14  minReplicas: 0
15  maxReplicas: 8
16  volumes:
17    - name: renv-cache
18      persistentVolumeClaim:
19        claimName: renv-cache
20  volumeMounts:
21    - mountPath: /renv
22      name: renv-cache
```

- Ray allows us to define worker nodes with specific capabilities
- Normally intended for nodes with GPUs available rather than CPU. But we use it to add different programming languages!
- Workers are scaled automatically as tasks come in, and the same worker is reused until all tasks are done. Pay the overhead only once!
- Workers are based off the usual images + the ray image.

```
macrobond_ray.py
1 from ray_provider.operators import SubmitRayJob
2
3
4 with dag:
5     source_macrobond = SubmitRayJob(
6         task_id="source_macrobond",
7         conn_id="kuberay",
8         entrypoint="Rscript /run_task.R kofpkg1 source_macrobond",
9         runtime_env={
10             "env_vars": {
11                 "https_proxy": "http://proxy.ethz.ch:3128",
12                 "DB_HOST": "db.kof.ethz.ch",
13                 "DB_PASSWORD": "{{var.value.prod_pg_password}}",
14                 "MACROBOND_API_PASSWORD": "{{var.value.mb_api_password}}",
15             }
16         },
17         num_cpus=0,
18         memory=0,
19         resources={
20             "r": 1,
21         },
22         outlets=[Dataset("macrobond")],
23     )
```

- The Ray provider makes the integration with Airflow very simple to achieve
- By defining our own entrypoint, we're able to execute arbitrary things in whatever programming language => R!
- Airflow still tracks task execution and status.

```
ray_helper.py
1 import ray
2
3
4 def ray_remote_r_task(*, env_vars=None):
5     if env_vars is None:
6         env_vars = {}
7
8     @ray.remote(resources={"r": 1}, env_vars=env_vars)
9     def run_r_task(package, task):
10         import rpy2.robjects as robjects
11         from rpy2.robjects.packages import importr
12
13         r = robjects.r
14         r.options(**{"renv.consent": True})
15         base = importr("base")
16         renv = importr("renv")
17
18         renv.install(package, prompt=False)
19         tasks_dir = base.system_file("tasks", package=package)
20
21         script = base.paste0(tasks_dir, "/", task, ".R")
22         print(f"Executing script: {script}")
23         base.sys_source(script, envir=r.topenv())
24
25     return run_r_task
```

- It is possible to go through the python sdk to execute R code
- However: who wants to write R code in python?
- Due to some technicality, this code is actually not reusable and must be copy/pasted in each Airflow task

Ray: Pros and Cons (so far)

Pros

- Automatic scaling of workers based on the amount of tasks
- Reduces overhead of starting/stopping containers
- Great capabilities to further distribute ML/AI models
- Flexible enough to allow arbitrary code execution – can be used with any language

Cons and open challenges

- Logs are not in Airflow anymore, need to be retrieved via specific logging infrastructure
- Keeping worker images up-to-date can be a challenge
- With great power, comes great complexity

Questions / Demo



ETH Zürich
KOF Institut

Leonhardstrasse 21
8092 Zürich
www.kof.ethz.ch

Insights empowering decisions