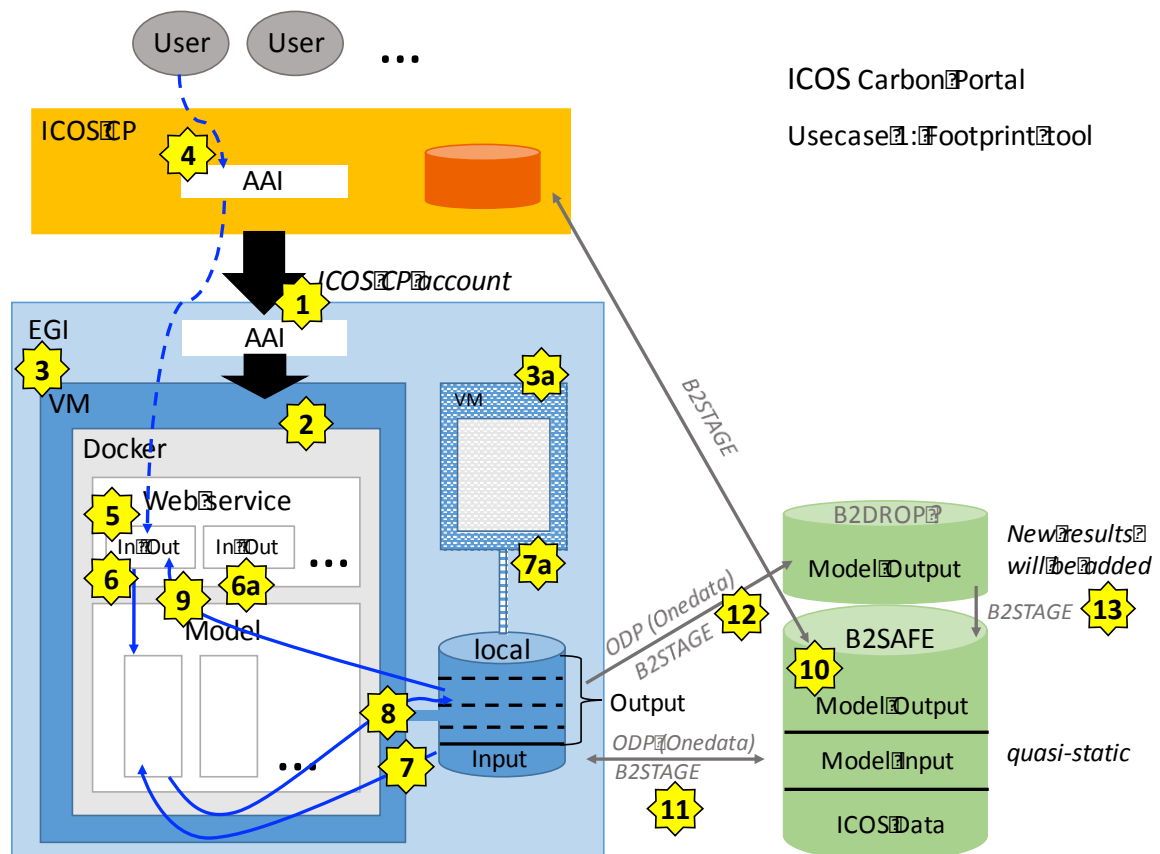


## ICOS Carbon Portal Usecase1: Footprint Tool

This is a short description of the work- and dataflow, storage and computation requirements of the usecase and its implementation with EGI and EUDAT services. It should only serve as a basis for discussions between EGI, EUDAT and ICOS on the implementation of the usecase. The strategy is still open for discussion and the description is not yet complete. **More information will be added as required during the discussion.**



### General work- and dataflow:

(Numbers correspond to those in figure above)

(Open questions in italic letters)

Workflow for setting up and providing the footprint tool service:  
(thick black arrows)

1. ICOS CP instantiates a VM with attached block storage in the EGI FedCloud. A robot certificate associated to ICOS CP is used for the authentication and authorization in the EGI FedCloud.
2. ICOS CP runs the web service and the model computations in a docker on the VM.
3. The VM is running continuously and has a 'permanent' IP address as long as it exists. This IP address is registered in the DNS of ICOS CP and the service is hosted on a subdomain of icos-ip.eu

**3a.** In case several VMs are required to provide the service, *the load balancing is done either at ICOS CP or where and how?*

Workflow for a user of the footprint tool service:  
(blue dashed arrows)

**4.** The user accesses the service at ICOS CP. Authentication and authorization of the user is handled at ICOS CP (HTTP cookie-based authentication). ICOS CP takes care of all interactions with the EGI FedCloud on behalf of the user.

**5.** The user is directed to the web service running in the docker on the VM.

Dataflow in the footprint tool:  
(blue arrows)

**6.** The web service provides the user interface to start model runs and visualize results. Input parameters are selected on the web page and passed to the script that starts the model run in the same docker.

**6a.** The web service will serve several users in parallel and launch separate specific model runs.

**6b.** In an intermediated step a data base (run either in the same or an additional docker on the same VM) is accessed to search a catalogue of existing model output in order to avoid recomputing already existing footprints and/or time series. *Further specifications are needed here.*

**7.** The model reads input files from the local storage directly attached to the VM (or directly from B2SAFE using ODP/onedata).

**7a.** In case of several VMs the input files need to be either stored locally at the individual VM or in a storage accessible from all VMs (based on ODP).

**8.** Model output is written to a run-specific directory in the local storage directly attached to the VM or (8a) to a common storage as in 7a.

**9.** Model output is displayed on the web page and the user can also download the results to her/his local computer.

Data storage and transfer:  
(grey arrows)

**10.** Quasi-static input data (incl. metadata) are stored on the ICOS CP server and for longterm storage in B2SAFE (data transfer using B2STAGE).

**11.** A copy of the input data and already available output data is kept on the storage attached to the VM or provided via OPD.

**12.** The new model output is added to (or merged with) already available model output in a storage for fast sharing that is globally accessible, independently from the specific VM (e.g. B2DROP or?), transfer using ODP or via B2SHARE.

**13.** Model output data (incl. metadata) are transferred regularly to B2SAFE (using B2STAGE?) and archived in B2SAFE for longterm storage.

*The strategy to attach a PID (or DOI) to the model output (and user-specified request) is not yet included.*

### **Workflow inside the footprint model:**

#### Scenario 1: ICOS sites

- Precomputed particle location files exist
- Footprints and time series exist for most sites and time ranges and emission types
- Check availability and compute only missing parts
- Display results

#### Scenario 2: New sites

- No precomputed particle location files exist
- Full STILT run required (compute particle location, footprint, time series)
- Display results

#### Detailed workflow:

- Select time range and station (name-id or latitude/longitude)
- Check availability of time series (csv-file) for full time range and all emission types
  - For each date (year, month, day, hour):
    - Check availability of aggregated footprints (netCDF-files)
    - If not available
      - [Compute particle location file] for scenario 2 only
      - Compute footprint for this date and store (in netCDF-file)
      - Compute concentration for this date based on footprint and emissions, append result to csv-file

-> Display time series and animation of aggregated footprints in web service

## Storage requirements:

- Input datasets (quasi-static, updated every 6-12 months)
  - Emissions (EDGAR, VPRM and more)
  - Initial and boundary concentration data
  - Meteorology
  - Particle location files (precomputed for many stations and dates)
  - Observation time series

Input dataset		no of files <u>per year</u>	File size	Storage <u>per year</u>	Type
Emissions	EDGARv4.1	3	6.3 GB	19 GB	netCDF
	EDGARv4.3	250	1-10 MB	< 1 GB	netCDF
	VPRM	10	1-6MB	< 1 GB	netCDF
	others	3-10	20-100 MB	< 1 GB	netCDF
Boundary		5	1-10 GB	20 GB	netCDF
Meteo		12	17G	205 GB	arl (binary)
Particle location	per station	1-3 hourly	2920-8760	0.5-5 MB	R-object
	all available	59 stations	~ 250000	0.5-5 MB	

Table shows estimated storage requirements for 1 year, but computation of more years will be added.

- Output datasets  
(model runs access output data from previous runs and eventually add new files)
  - Aggregated Footprints
  - Concentration time series
  - Particle location files for new sites (produced in a full STILT run)

Output dataset		no of files <u>per year</u>	File size	Storage <u>per year</u>	Type
Footprints	per station	1-3 hourly	2920-8760	1.5 MB	netCDF
	all available	59 stations	~ 250000	1.5 MB	
	user requests	new sites			
Time series	per station	1-3 hourly	1-52	1-50 MB	csv (ascii)
	all available	59 stations	~ 100	1-50 MB	
	user requests	new sites			
Particle location	per station	1-3 hourly	2920-8760	0.5-5 MB	R-object
	user requests	new sites			

Table shows estimated storage requirements for 1 year, but computation of more years will be added. Furthermore, user requests will initiate computation of additional footprints and time series for existing and new sites.

### **Computational requirements:**

Test on a linux cluster at MPI for Biogeochemistry in Jena and VM in EGI FedCloud  
*(will add specification here)*

Full STILT run:

3 GB memory per job

670 CPUs per footprint

1700 CPUh per station per year

*(more detailed information on cpu time etc. will be added)*

Model runs for individual stations and users requests are separate jobs therefore parallel processing is possible and required for better performance.