# Big Data, organisation and analysis
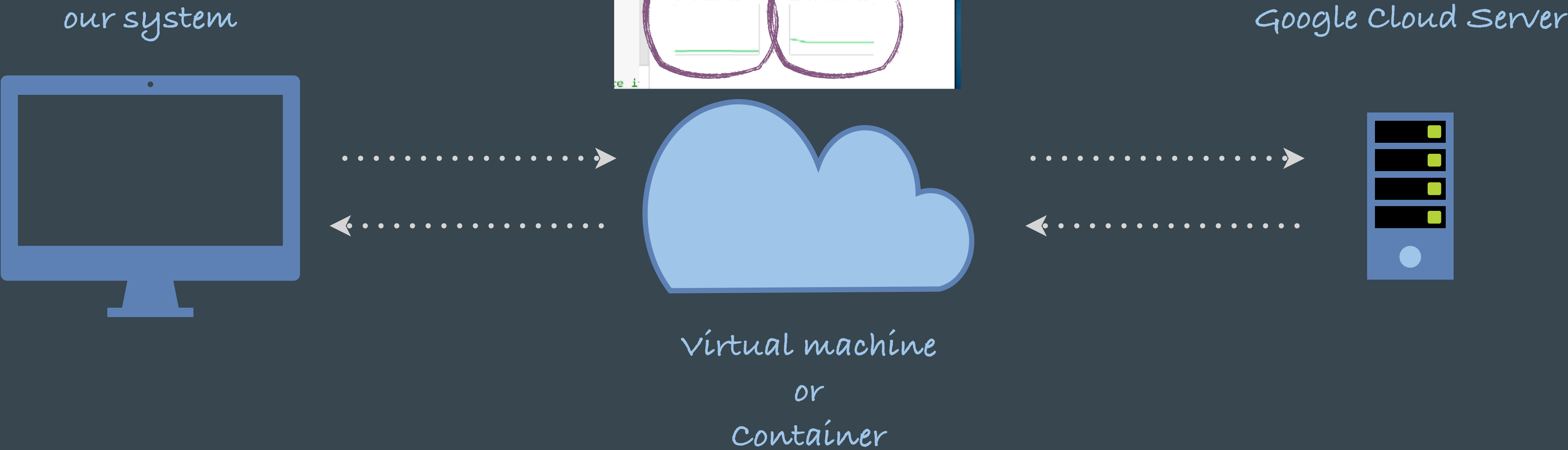
Building a REST API service

Steffen M. Noe, Spring 2023

# Recap on what we learned

- So far, we have been learning...

  - What is virtualisation and the link with VMs and Containers

  - How these are forming Cloud computing services (e.g., Colab Notebooks)

  - How to create a "middleware" to access data

  - How to access REST APIs

# Google Cloud service

## Principle



our system

Google Cloud Server

Virtual machine
or
Container

# Data access via API

Principle



Google Cloud virtual service

Client system

Kaggle data server

API request

API response

# REST API access via requests package

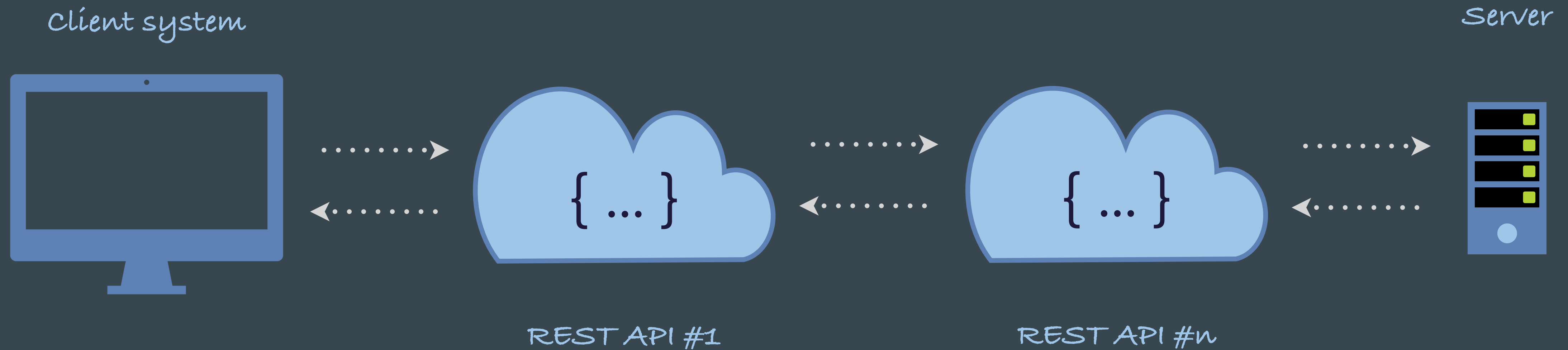## Principle



Client system

Server

REST API

API request

API response

writing data to file on the virtual server

# REST APIs can be chained

There can be more than one API layer, it can be many!

Client system

Server

{ ... }

{ ... }

REST API #1

REST API #n

# We need to know also how to build the "other side" - the server

- REST API is a client - server system

- So far we know how to use it as a client

Client system

Server

{ ... }

REST API

# Build a simple Docker server with REST API

- Step 1: Download Docker Desktop

- Step 2: Choose a Docker image (Some Linux version as example)

- Step 3: Create a simple REST API server

# REST API

Principle

Client system

Server

{ ... }

REST API