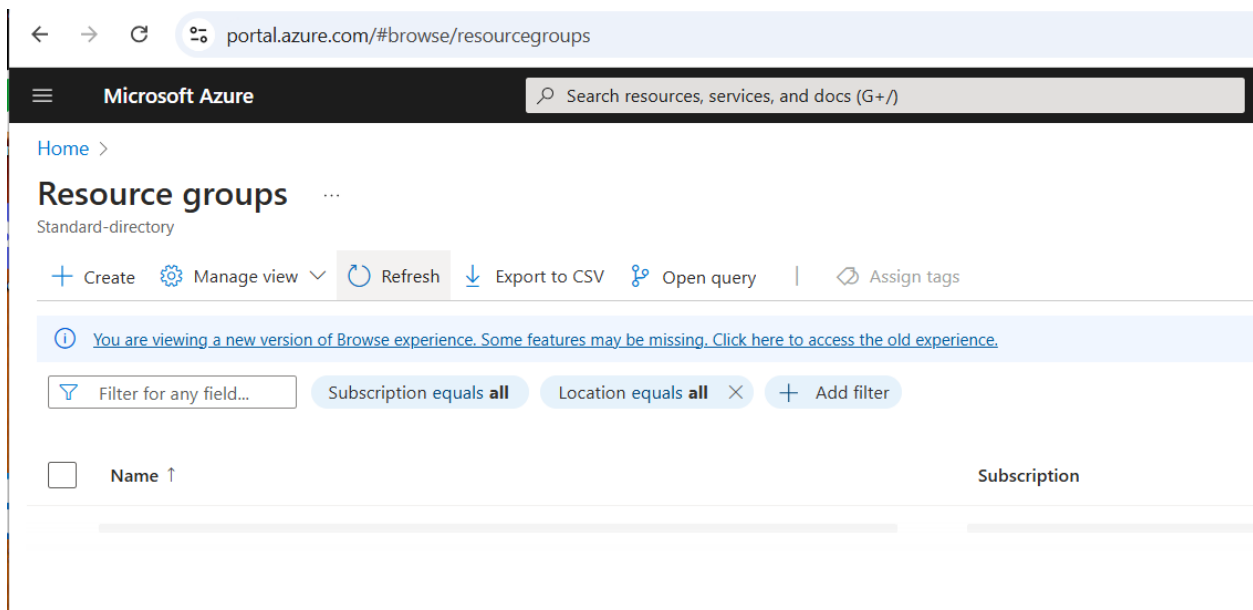# Modular RAGs

Finn Vilsbæk

fvil@eadania.dk

1

# Subjects for today

- Azure AI Services setup
- Running the example code base

- Prerequisites: An active Azure account, preferably with a pay-as-you-go subscription and your own dedicated ressource group.

- Get the files you will need for this workshop from here: http://panmedia.dk/en-US/rag-workshop

  ..and unzip them into a local folder on your disk.
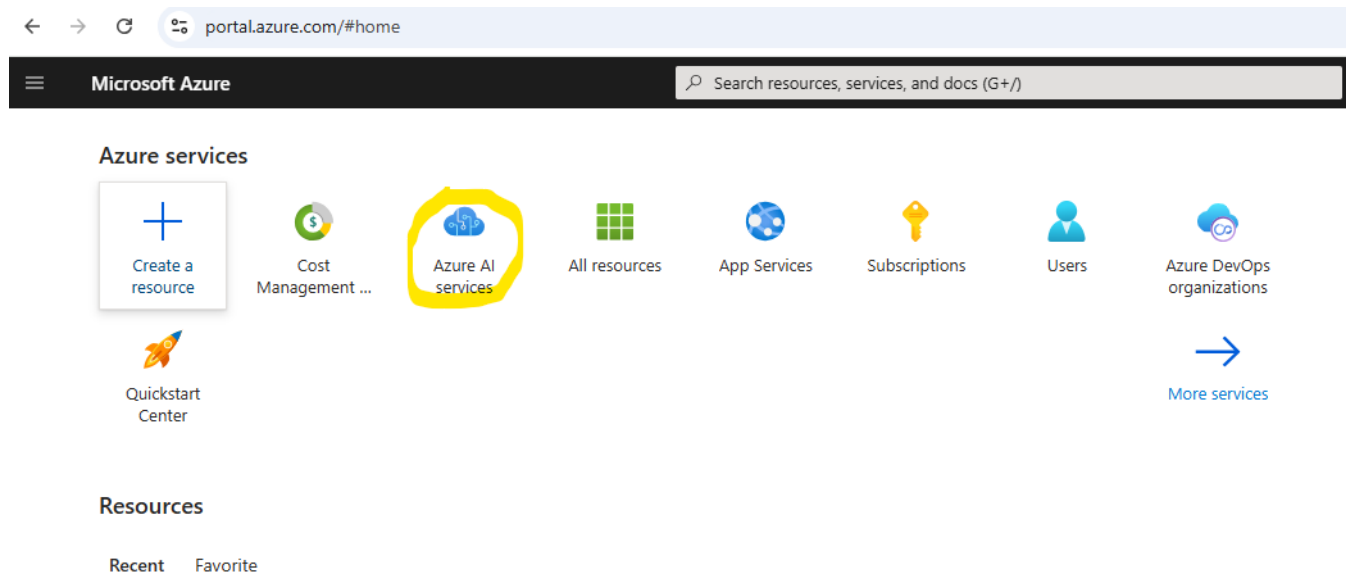
# Modular RAGs

- In order to get the example code up and running on your local machine, you first need to have an active Azure Account – your student account with your institution should work fine. Create a resource group first, as you need to have a place to store your services in. You can create this by searching for 'resource group' in the search bar.

# Modular RAGs

- From your Azure Portal Home Page, pick Azure AI Services.
  Link: https://portal.azure.com/#home

# Modular RAGs

- Bob is your uncle today, because AI Search and Computer Vision, the two main services we need are pretty much at the top of the list.

# Modular RAGs

- Create a Computer Vision instance.

  Be sure to pick the free F0 tier.

# Modular RAGs

- The Azure AI Search service is a little trickier to set up. Change the price tier via the link at the bottom of the page from standard to..

# Modular RAGs

- .. free, which offers you a whopping three indexes – more than we need!

# Modular RAGs

- Setting up a search index is your next task. Simply go to your AI Search instance when the resource is created, and click on 'Add Index' – choose the normal index, not the JSON variant.

# Modular RAGs

- Give your Index a memorable name, and click on 'Add Field'.

# Modular RAGs

- Create an url field. Check all the boxes, so the field gets a standard Lucene analyzer.

**Index Field** ✕

Field name *

url

Type ⓘ

Edm.String ⌄

**Configure attributes**

☑ Retrievable

☑ Filterable

☑ Sortable

☑ Facetable

☑ Searchable

Analyzer

Standard - Lucene ⌄

Save   Cancel

# Modular RAGs

- Create a field 'contentVector' of type Collection(Edm.Single) with the settings shown here, and also create a vector search profile.

**Index Field** ✕

Field name *

contentVector

Type ⓘ

Collection(Edm.Single) ⌄

**Configure attributes**

☑ Include in storage ⓘ

☑ Retrievable

☑ Searchable

Dimensions * ⓘ

1024

Vector search profile *

No vector search profiles

Create

Save  Cancel

# Modular RAGs

- You will also need to create an algorithm configuration for your Vector Search Profile.

# Modular RAGs

- The Vector algorithm should have the settings shown here:
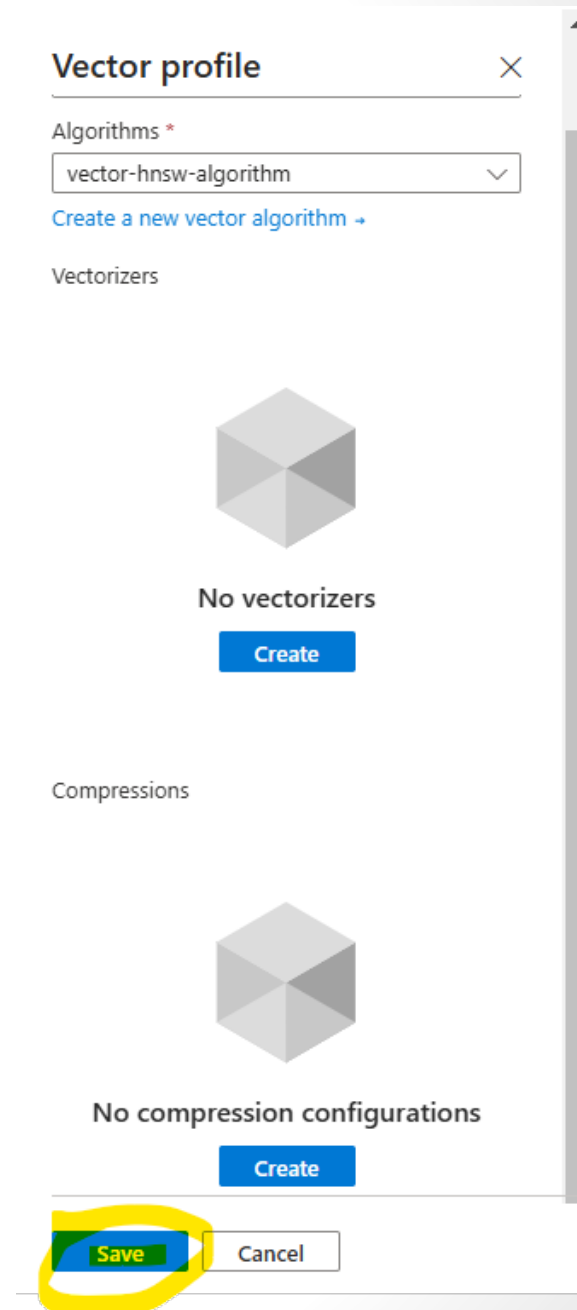
# Modular RAGs

- Now you need to save your new Vector profile – make sure that your algorithm is selected in the dropdown box and click Save.

  You don't need to create a Vectorizer or a Compression configuration.

# Modular RAGs

- Now the system will return to the previous screen, and you can save the contentVector index field, along with the Vector profile you have just created. NB: **make sure** that the boxes are ticked here!

**Index Field** ✕

Field name *

contentVector

Type ⓘ

Collection(Edm.Single) ⌄

**Configure attributes**

☑ Include in storage ⓘ

☑ Retrievable

☑ Searchable

Dimensions * ⓘ

1024

Vector profile * ⓘ

vectorSearchProfile ⌄

Create new vector profile →

Save    Cancel

# Modular RAGs

- Lastly, you need to create the index. If you have made the index precisely as indicated in the previous screendumps, the boxes are ticked as shown here, and you can hit 'Create'. **NB: if you click any of the checkboxes in this view, you have to make the index from the starting point again** – it's a bug in the Azure interface, sorry guys ☹

# Modular RAGs

- Now, you can set up your unique API keys and endpoints from your own Azure account in the two solution projects. The appsettings.json file in the Console project:
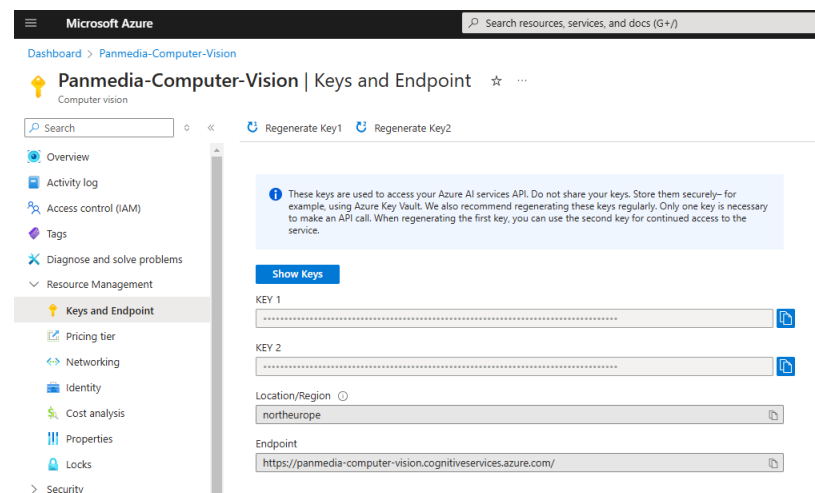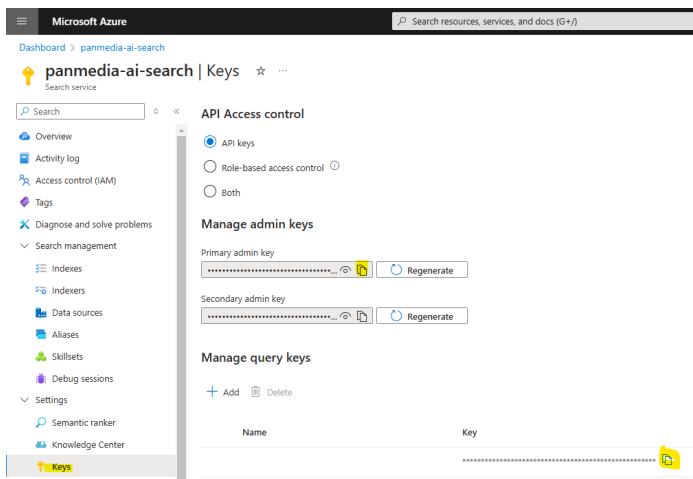


.. and the appsettings.json file in the WebApplication project:

# Modular RAGs

You can copy the keys and endpoints you need from each resource from Settings >> Keys and Resource Management >> Keys and Endpoints. Click on the icon marked with yellow to copy the key to the clipboard.

# Modular RAGs

NB: the search endpoint in the console app indexer is special:

"AzureAiSearchEndpoint":

https://YOUR_NAMED_SEARCH_INSTANCE.search.windows.net/indexes/YOUR_NAMED_INDEX/docs/index?api-version=2023-11-01

Here, you can set up your own names instead of the capital letters.

# Modular RAGs

- Important: set your primary admin key for ai search in the Console project, and use the more lowly query key in the Web project. If you don't, the Console App will not index correctly.
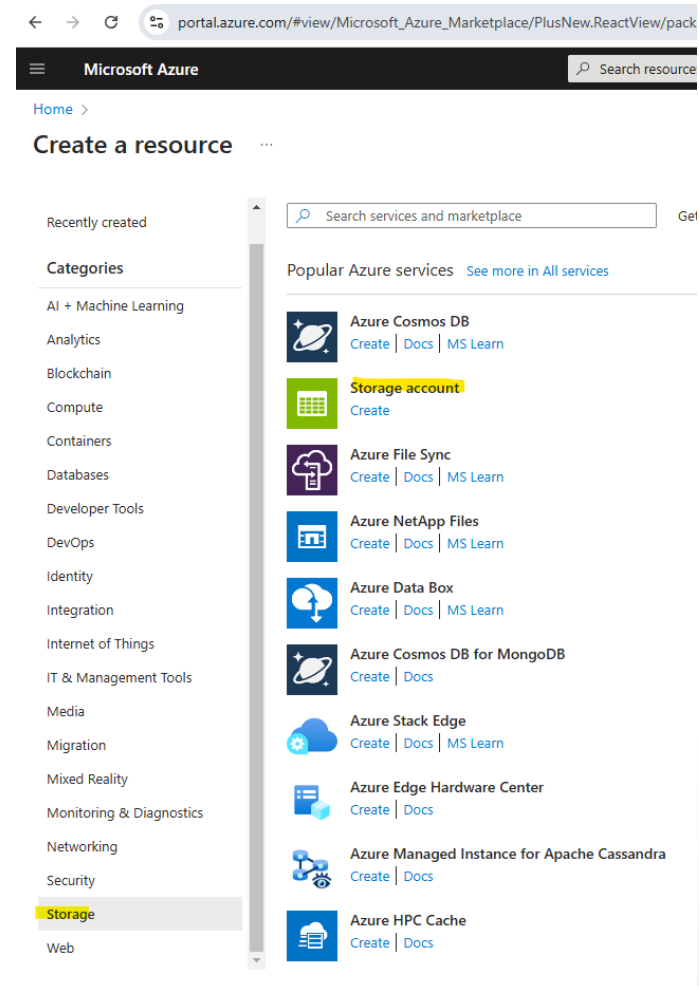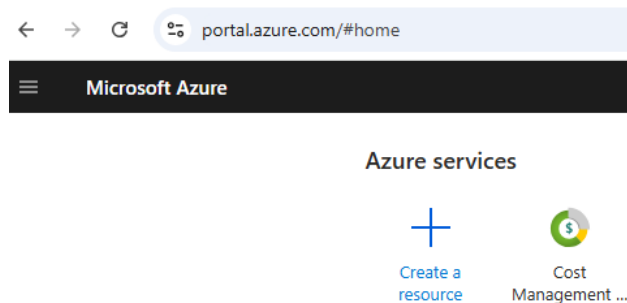
# Modular RAGs

- In the file 'VectorRepository.cs' in the Web project, you can play around with the value for k in line 28 to 4 or 5 instead of 3. This will later bring up three, four or five pictures from the urls in the search index. There aren't many images though in the sample folder, and with the value of 3 you will see less 'hallucinatory' pictures at the bottom of the list ☺

# Modular RAGs

- Now for the pretty pictures part. You will need to set up a small blob storage on Azure, since the indexing we need to do is a lot easier on that platform. Go to https://portal.azure.com/#home, and click on 'Create a resource'. Then choose 'Storage Account'.





23

# Modular RAGs

If you are unable to create a Blob storage container on your standard Azure student account – I know this can be pretty troublesome - then you can use my Blob storage account in the appsettings.json file:

```
       AzureOpenAiService.cs                                          appsettings.json
https://json.schemastore.org/appsettings.json
 8          "AzureBlobContainerUrl": "https://panmediablob.blob.core.windows.net/images"
 9        }
10
```

If you want to try to create your own Blob storage container, you can proceed with the next six slides.

# Modular RAGs

- I chose the cheapest possible options. Don't worry too much here, it most probably isn't going to break the bank account, but it is necessary in order to make the App work.

# Modular RAGs

- Check your settings – make sure that your storage region is relevant to the region that your other resources belong in, and make sure that your primary service is set to include 'Azure Blob Storage'.

# Modular RAGs

- Next, we will need a container for our images. Click on the plus sign next to 'Container', and select the most permissive anonymous access level, since we want our Console App to be able to see and enumerate the container contents.

# Modular RAGs

- Take the images from the zip file you have downloaded - and upload them into your container. There should be 60 images in all.

# Modular RAGs

- Now, with the keys and endpoints correctly set up in both instances of appsettings.json, we should be able to run the indexer from the console. The project you have downloaded will need to be set to 'ConsoleAppIndexer' in the dropdown in VS 2022, and you can click on the green triangle to run the indexer.



- If it doesen't work as expected, get my attention and I will take a look at your setup.

# Modular RAGs

- If everything is working, you should see the indexer crawling the Azure Blob store image container:

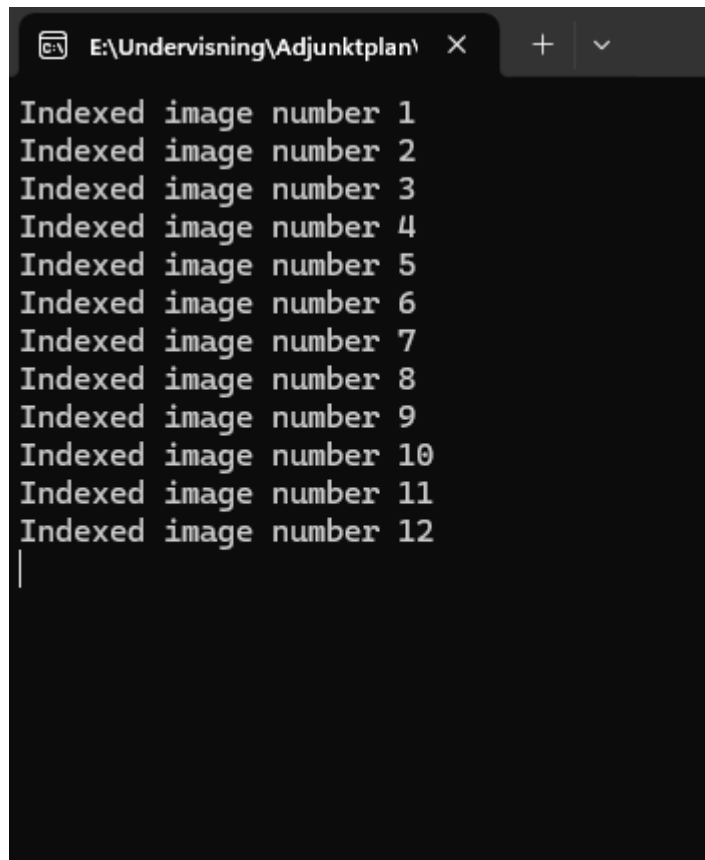There is a delay of three seconds between each index entry produced, because the free tiers of the services we are using have a limit of 20 requests per minute. Let the indexer do its job in the foreground and leave your machine to it until 60 images have been processed.

```
E:\Undervisning\Adjunktplan\    ×    +    ∨

Indexed image number 1
Indexed image number 2
Indexed image number 3
Indexed image number 4
Indexed image number 5
Indexed image number 6
Indexed image number 7
Indexed image number 8
Indexed image number 9
Indexed image number 10
Indexed image number 11
Indexed image number 12
```

# Modular RAGs

- Once the indexer is done, choose the Web Application from the top menu, and click on the green triangle to start the Web App.



- Click on 'Search' to enter a search term for an animal.

# Modular RAGs

- Here, I have searched for cats:

# Modular RAGs

- .. and here I have searched for crabs. The vector search will return the highest probability of a match first (the only picture of a crab in the image collection) and it will 'hallucinate' the closest match it thinks it can find for the next two images.

# Modular RAG's

- Learn more about Azure AI Services

- Microsoft Learn AI Services landing page: https://learn.microsoft.com/en-us/azure/ai-services/

- AI Search: https://learn.microsoft.com/en-us/azure/search/

- Computer Vision: https://learn.microsoft.com/en-us/azure/ai-services/computer-vision/