

Annotation methodology

What is in this file?

This document contains a set of instructions pertaining the annotation of the attached data.

Context

This data is a collection of queries and images based on the MS COCO dataset. The dataset contains ~ 100k images and text descriptions of the images. We have selected a set of multiple texts for which we have determined a relevant image set. In order to perform our study we need to annotate the set of returned images obtained through our search engine.

How the data is structured?

Alongside this document you will find multiple folders containing images.

Each folder name respects the schema: <QUERY_ID>_<QUERY_TEXT>.

The <QUERY_ID> represents the unique identifier of the query. These are values between 1 and 10000. You might find attached only a subset of them. The <QUERY_TEXT> represents the text specific to the query.

Inside each folder you will find the set of returned images for the specific query.

The naming schema of the images is the following <IMAGE_RETRIEVAL_INDEX>_<IMAGE_SCORE>.jpg.

The <IMAGE_RETRIEVAL_INDEX> represents the associated index in our internal search engine. You might see that some values are skipped altogether. This is because our search engine is based on text-to-text search and some images contain multiple texts. When ranking, the same image might have appeared multiple times. Which is why, some of the positions are skipped as they have already been considered.

The <IMAGE_SCORE> represents the score given by the retrieval engine.

How do we annotate data?

PLEASE READ THE FULL INSTRUCTIONS BEFORE STARTING:

In order to annotate a query we need to:

1. Decide the split of data we are going to cover and discuss it with the team as to not cover the same data twice. Ex: "I will annotate all queries from the id 5000 to 6000". After confirming it with the team we can start our annotation process.
2. You will find attached a spreadsheet with the following columns: query_id, last_id. You will need to open this while annotating as you will need to save some information in it.

3. Select the query folder you want to annotate. Enter the corresponding folder. The <QUERY_TEXT> contains the items you will be searching for:

Ex: "A floor looks rounded in the neat clean living room".

3. Given the previous QUERY_TEXT we will now go through each of the images in the folder. We will DELETE all the images which do not contain exactly what the query has. We will note down the <IMAGE_RETRIEVAL_INDEX> of the last match we have observed.

Obs.

The first image is most probably correct as it is the ground-truth for that query so you should have at least 1 match.

4. Note down in the spreadsheet the <QUERY_ID> and the <IMAGE_RETRIEVAL_INDEX> of the last identified match. (please do them in order :)).

5. After finishing this entire process, you will create an archive of the annotated queries, and the spreadsheet and send them over.

Suggestions:

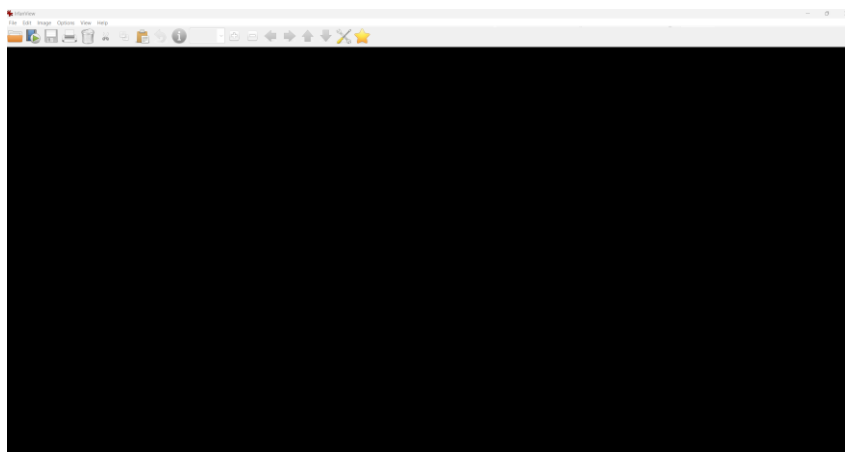
Install IrfanView. This is a image viewer/editor/organizer/convert available for Windows <https://www.irfanview.com/64bit.htm>

This will allow you to make use of IrfanView Thumbnail and allow you to speed up your annotation process by a lot!

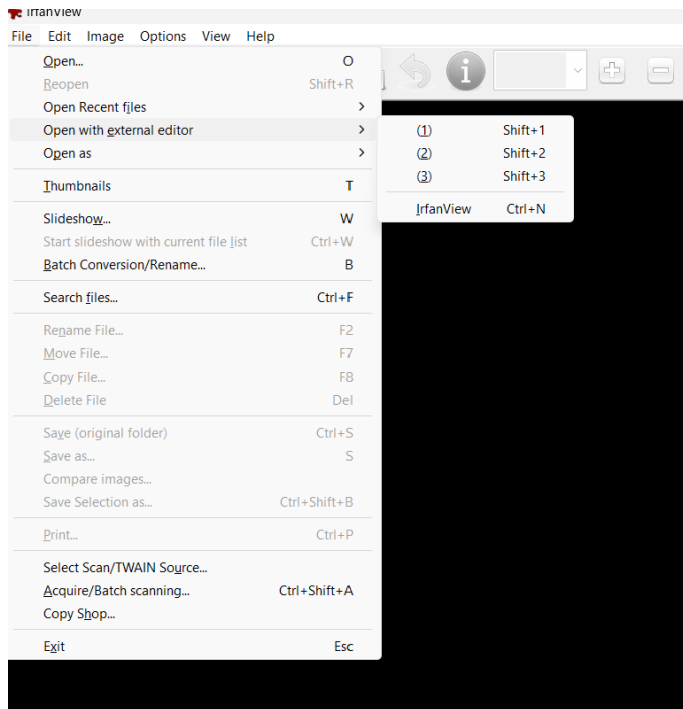
Setup steps for IrfanView:

1. Install Irfanview from the previous link.

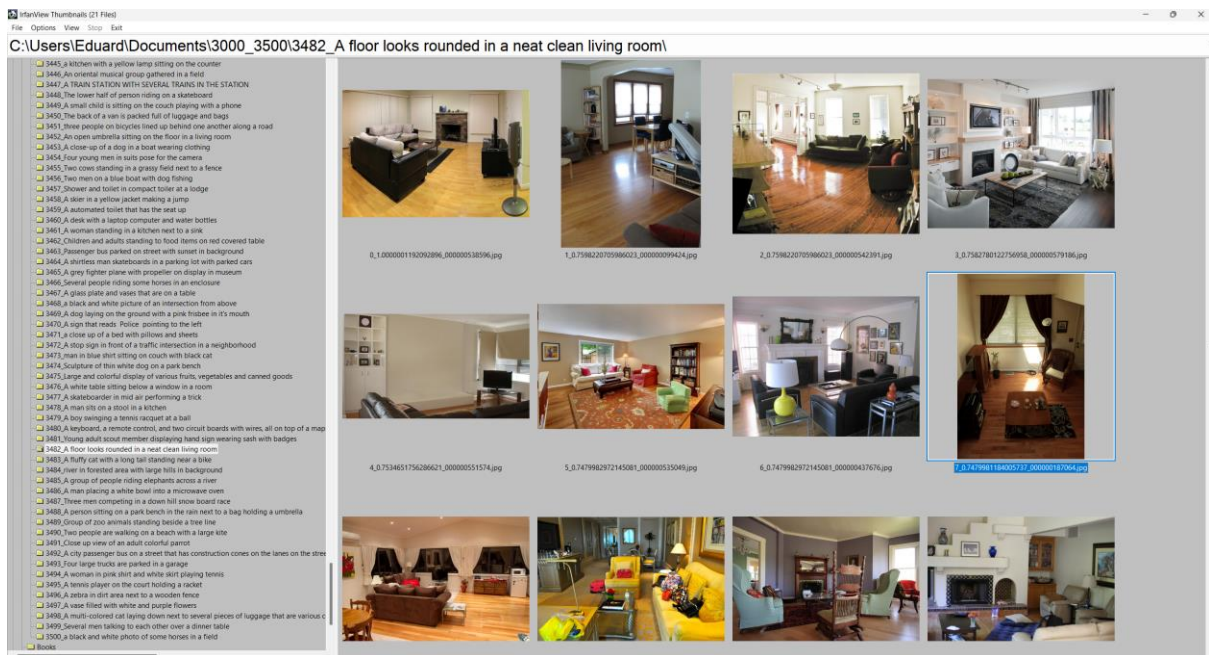
2. OPEN IRFANVIEW



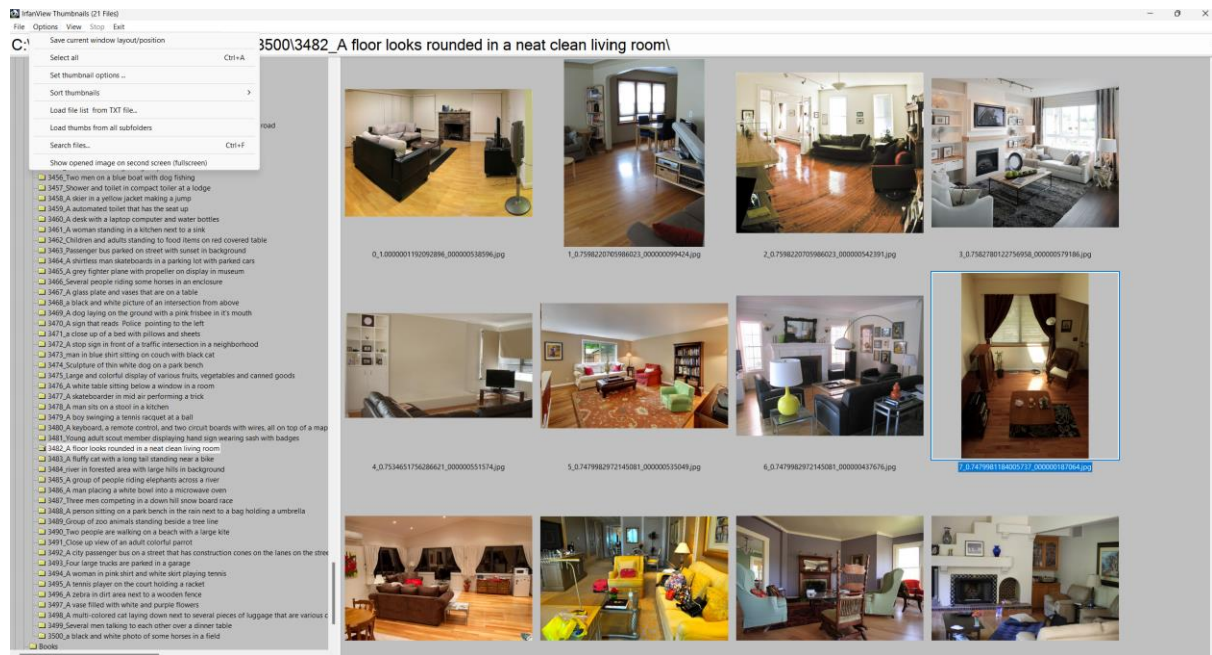
2. GO TO FILES>THUMBNAILS.



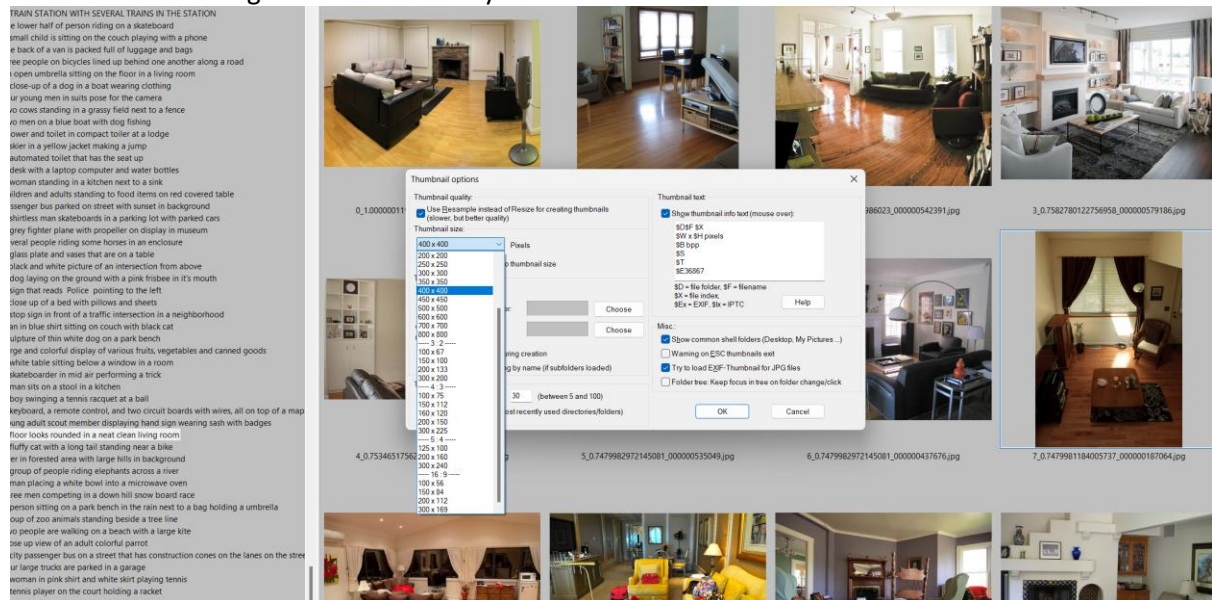
3. Go to the folder you want to annotate.



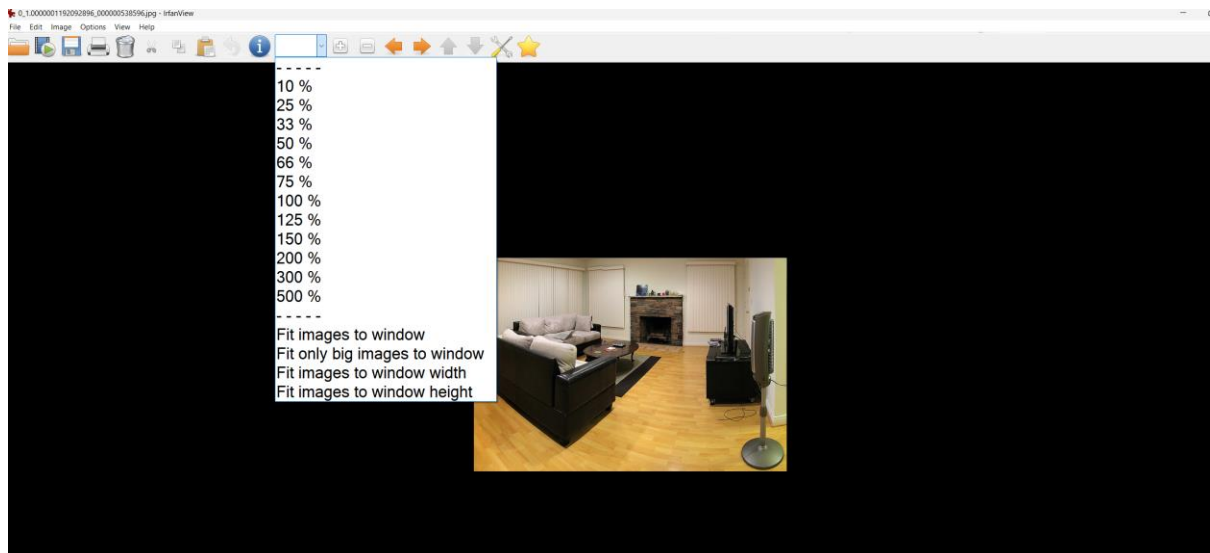
4. Set the thumbnail size to a more appropriate size in order to render multiple on the same page. To do this go to Options> Set Thumbnail Options...



5. Set the right size for your screen from the thumbnail size selector

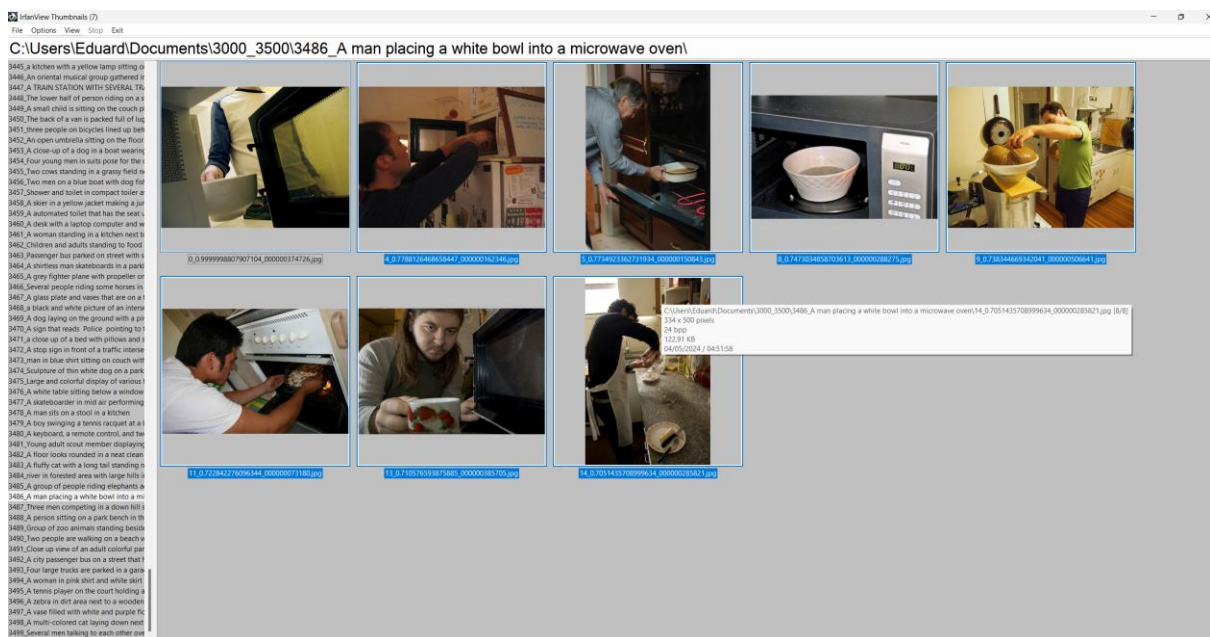


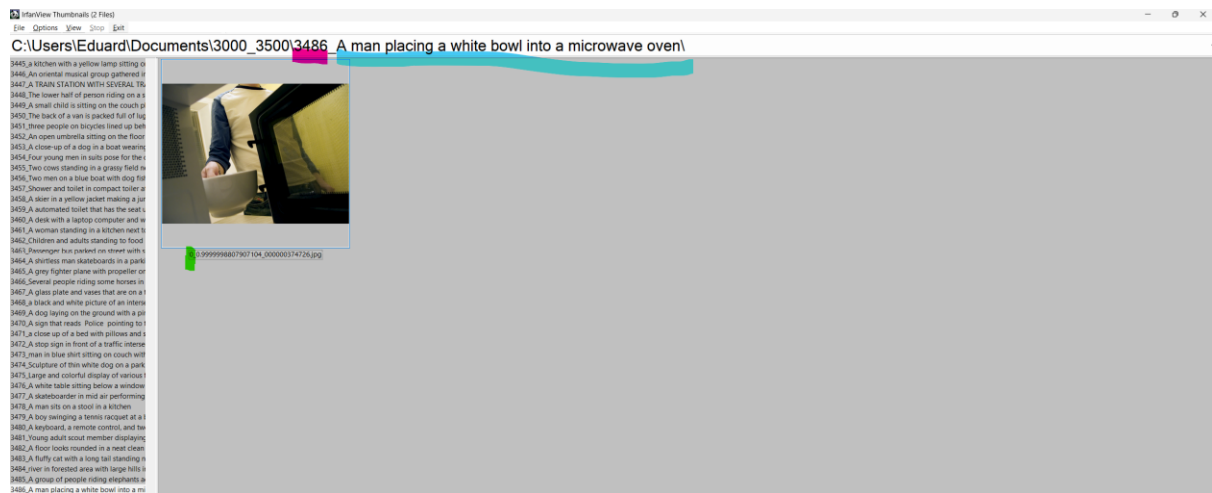
6. If you want to be able to zoom in even more on an image you can double click on it to open the image viewer. You can afterwards set the zoom level in the top bar



To keep this zoom for all the images you double click on go to View > Lock Zoom.

7. Select the images you want to delete and delete them from irfanview with DEL key





8. Open up the spreadsheet and go the the Query ID

	A	B	C	D
3475	3473	0		
3476	3474	1		
3477	3475	0		
3478	3476	56		
3479	3477	2115		
3480	3478	7		
3481	3479	338		
3482	3480	0		
3483	3481	0		
3484	3482	21		
3485	3483	37		
3486	3484	7		
3487	3485			
3488	3486			
3489	3487			
3490	3488			
3491	3489			
3492	3490			
3493	3491			
3494	3492			
3495	3493			
3496	3494			
3497	3495			
3498	3496			
3499	3497			
3500	3498			
3501	3499			

9. Write the last identified matching IMAGE_RETRIEVAL_INDEX.

manual_gt

Fişier Editează Afişează Inserează Formatează Date Instrumente Extensii Ajutor

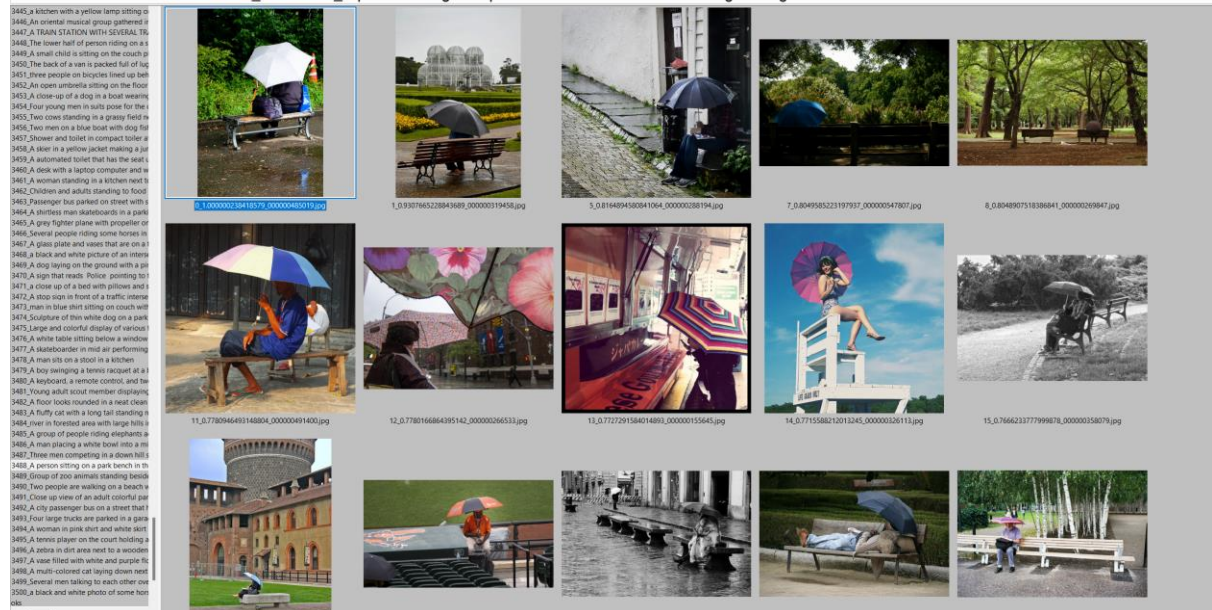
100% | lei % 123 | Presta... | 10 | B I A

D3485

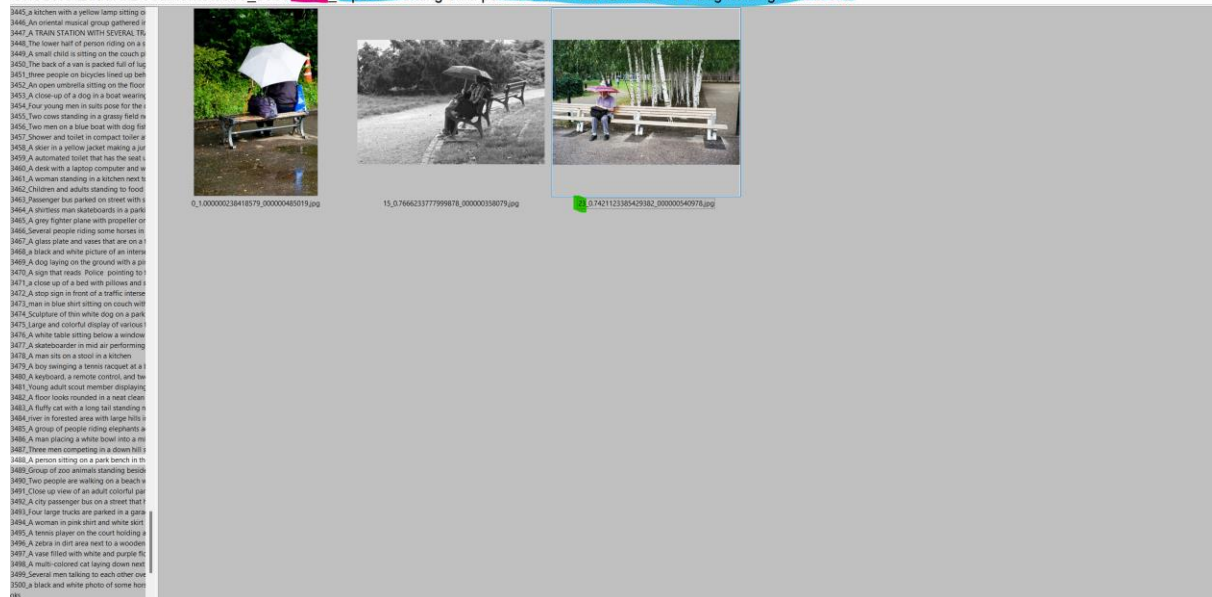
	A	B	C	D
3475	3473	0		
3476	3474	1		
3477	3475	0		
3478	3476	56		
3479	3477	2115		
3480	3478	7		
3481	3479	338		
3482	3480	0		
3483	3481	0		
3484	3482	21		
3485	3483	37		
3486	3484	7		
3487	3485			
3488	3486	0		
3489	3487			
3490	3488			
3491	3489			
3492	3490			
3493	3491			
3494	3492			
3495	3493			
3496	3494			
3497	3495			
3498	3496			
3499	3497			
3500	3498			
3501	3499			

Example 2

C:\Users\Eduard\Documents\3000_3500\3488_A person sitting on a park bench in the rain next to a bag holding a umbrella\



C:\Users\Eduard\Documents\3000_3500\3488_A person sitting on a park bench in the rain next to a bag holding a umbrella\



	A	B	C	D
3479	3477	2115		
3480	3478	7		
3481	3479	338		
3482	3480	0		
3483	3481	0		
3484	3482	21		
3485	3483	37		
3486	3484	7		
3487	3485	56		
3488	3486	0		
3489	3487	0		
3490	3488	2		
3491	3489			
3492	3490			
3493	3491			
3494	3492			
3495	3493			
3496	3494			
3497	3495			
3498	3496			
3499	3497			
3500	3498			
3501	3499			
3502	3500			
3503	3501			
3504	3502			
3505	3503			