

## JSONL Format

To finetune polyGemma we need dataset in  
JsonL format

JsonL is similar to Json, but each line is  
an independent Json Object

[  
  { ... }  
  { ... }  
  { ... } ]  
Each Line is a separate JSON



{  
  image : path to image  
  prefix : prompt  
  suffix : Response  
}

Each json contains 3 keys, image, prefix, suffix.

image: holds the path to associated image

prefix: prompt we are going to send to model.

suffix: Expected model response.

»

When used for Image Captioning

prefix could be describing the image &

suffix could be a man with a dog



  image : path to man holding dog

  prefix : describing the image

  suffix : A man holding a dog }

>> when used for Visual Question Answering (VQA)

Prefix could be what breed of dog is it

Suffix could be answer like a beagle

{

image : path to image

prefix : 'what breed of dog is it'

suffix : 'Beagle'

}

>> Object Detection :

Similar to other task, there is prefix & suffix.

For prefix we can have detect a dog, where  
key word detect is critical.

Suffix has specific format, 4 consecutive loc  
tags, followed by desired class name.

If there are multiple objects, another section  
appears, separated by semicolon

{

image : path to image

prefix : detect a dog

suffix : <loc 8> <loc 9> <loc 10> <loc 12>

dog

}

If there are multiple object detected, they are  
displayed after semicolon

{ image : path to image  
 prefix : detect a dog  
 suffix : <loc-><loc-><loc-><loc-> dog;  
           <loc-><loc-><loc-><loc-> dog

}

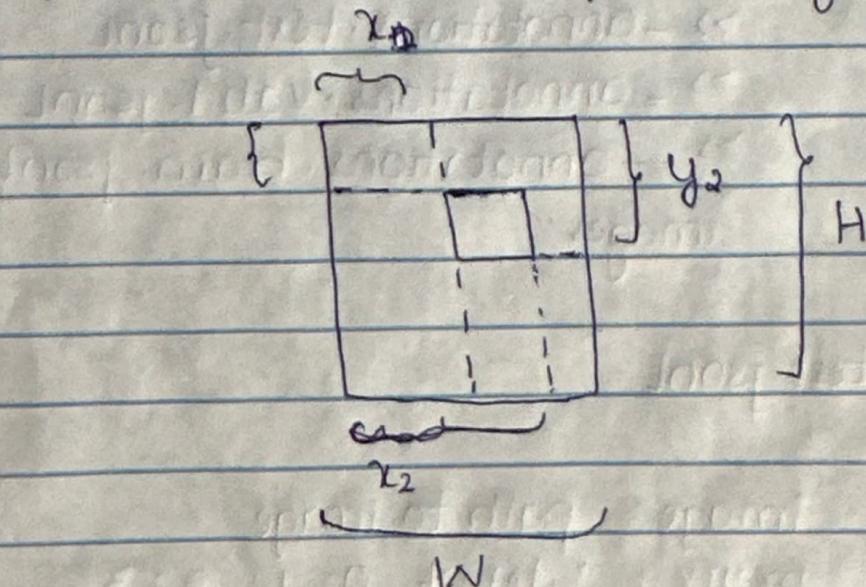
we can detect multiple object by listing them in suffix, separated by semicolon.

{ image : path to image.  
 prefix : detect a dog; cat; person  
 suffix : <loc-><loc-><loc-><loc-> dog;  
           <loc-><loc-><loc-><loc-> cat.

}

» The 4 loc tag, define the location of bounding box

Imagine a image, with detection defined by coordinates  $x_1, x_2, y_1, y_2$ , the image has height  $H$ , width  $W$



The first thing is to set coordinate in right order  
 $y_1, x_1, y_2, x_2$

$\text{bbox} = [y_1 \ x_1 \ y_2 \ x_2]$

Normalize the image

$\text{bbox} = \text{bbox} / [H, W, H, W]$

$\text{bbox} = \text{bbox} * 1024$

Eg: we place the location o padded to 4 digit

$\text{bbox} = [300, 400, 500, 600]$

$\langle \text{loc} 0300 \rangle \langle \text{loc} 0400 \rangle \langle \text{loc} 0500 \rangle \langle \text{loc} 0600 \rangle$

Finally we append class name after loc tag.

$\langle \text{loc} 0300 \rangle \langle \text{loc} 0400 \rangle \langle \text{loc} 0500 \rangle \langle \text{loc} 0600 \rangle \text{ dog}$

>> Dataset From Roboflow has following format:

dataset

>> -annotations-test.jsonL

>> -annotations-valid.jsonL

>> -annotations-train.jsonL

images

Eg: test.jsonL

{ image: path to image

prefix: "detect 0; 1; 2; 3...."

suffix:  $\langle \text{loc} 0048 \rangle \langle \text{loc} 0083 \rangle \langle \text{loc} 0683 \rangle$

$\langle \text{loc} 0944 \rangle 9$