# A. Pretrained Models

We specify details about all the pretrained models used, as well as the code-generation large language model:

- **GLIP [32]**. We use the implementation from the official GitHub repository[3]. In our experiments we use the GLIP-L (large) version. In order to adapt to new versions of PyTorch, we had to modify the CUDA implementation of some functions, as the repository relies on old versions of PyTorch. We provide our updated version of GLIP in our code.

- **MiDaS [45]**. We use the implementation from PyTorch hub[4], and use the "DPT_Large" version.

- **BLIP-2 [31]**. We tried both the implementation from the official repository[5] and the Huggingface one[6], with little difference between the two, being the former slightly more performant and the latter faster. In both cases, we used the Flan-T5 XXL version.

- **X-VLM [67]**. We used the official implementation[7], specifically the version finetuned for retrieval on MSCOCO.

- **GPT-3 for `llm_query`**. The GPT-3 model we use for the LLM query function is the `text-davinci-003` one. We use the official OpenAI Python API[8].

- **Codex**. The GPT-3 model we use for code generation is the `code-davinci-002` one.

See the code for more detailed implementation details.

# B. API

We provide the full API next, in Listing 1:

```
 1  class ImagePatch:
 2      """A Python class containing a crop of an image centered around a particular object, as well as relevant information.
 3      Attributes
 4      ----------
 5      cropped_image : array_like
 6          An array-like of the cropped image taken from the original image.
 7      left : int
 8          An int describing the position of the left border of the crop's bounding box in the original image.
 9      lower : int
10          An int describing the position of the bottom border of the crop's bounding box in the original image.
11      right : int
12          An int describing the position of the right border of the crop's bounding box in the original image.
13      upper : int
14          An int describing the position of the top border of the crop's bounding box in the original image.
15
16      Methods
17      -------
18      find(object_name: str)->List[ImagePatch]
19          Returns a list of new ImagePatch objects containing crops of the image centered around any objects found in the
20          image matching the object_name.
21      exists(object_name: str)->bool
22          Returns True if the object specified by object_name is found in the image, and False otherwise.
23      verify_property(property: str)->bool
24          Returns True if the property is met, and False otherwise.
25      best_text_match(option_list: List[str], prefix: str)->str
26          Returns the string that best matches the image.
27      simple_query(question: str=None)->str
28          Returns the answer to a basic question asked about the image. If no question is provided, returns the answer
29          to "What is this?".
30      compute_depth()->float
31          Returns the median depth of the image crop.
32      crop(left: int, lower: int, right: int, upper: int)->ImagePatch
33          Returns a new ImagePatch object containing a crop of the image at the given coordinates.
34      """
35
36      def __init__(self, image, left: int=None, lower: int=None, right: int=None, upper: int=None):
37          """Initializes an ImagePatch object by cropping the image at the given coordinates and stores the coordinates as attributes.
```

---

[3]https://github.com/microsoft/GLIP
[4]https://pytorch.org/hub/intelisl_midas_v2/
[5]https://github.com/salesforce/LAVIS/tree/main/projects/blip2
[6]https://huggingface.co/Salesforce/blip2-flan-t5-xxl
[7]https://github.com/zengyan-97/X-VLM
[8]https://openai.com/blog/openai-api

```python
            If no coordinates are provided, the image is left unmodified, and the coordinates are set to the dimensions of the image.
            Parameters
            -------
            image : array_like
                An array-like of the original image.
            left : int
                An int describing the position of the left border of the crop's bounding box in the original image.
            lower : int
                An int describing the position of the bottom border of the crop's bounding box in the original image.
            right : int
                An int describing the position of the right border of the crop's bounding box in the original image.
            upper : int
                An int describing the position of the top border of the crop's bounding box in the original image.

            """
            if left is None and right is None and upper is None and lower is None:
                self.cropped_image = image
                self.left = 0
                self.lower = 0
                self.right = image.shape[2]  # width
                self.upper = image.shape[1]  # height
            else:
                self.cropped_image = image[:, lower:upper, left:right]
                self.left = left
                self.upper = upper
                self.right = right
                self.lower = lower

            self.width = self.cropped_image.shape[2]
            self.height = self.cropped_image.shape[1]

            self.horizontal_center = (self.left + self.right) / 2
            self.vertical_center = (self.lower + self.upper) / 2

    def find(self, object_name: str) -> List[ImagePatch]:
        """Returns a list of ImagePatch objects matching object_name contained in the crop if any are found.
        Otherwise, returns an empty list.
        Parameters
        ----------
        object_name : str
            the name of the object to be found

        Returns
        -------
        List[ImagePatch]
            a list of ImagePatch objects matching object_name contained in the crop

        Examples
        --------
        >>> # return the children
        >>> def execute_command(image) -> List[ImagePatch]:
        >>>     image_patch = ImagePatch(image)
        >>>     children = image_patch.find("child")
        >>>     return children
        """

    def exists(self, object_name: str) -> bool:
        """Returns True if the object specified by object_name is found in the image, and False otherwise.
        Parameters
        -------
        object_name : str
            A string describing the name of the object to be found in the image.

        Examples
        -------
        >>> # Are there both cakes and gummy bears in the photo?
        >>> def execute_command(image)->str:
        >>>     image_patch = ImagePatch(image)
        >>>     is_cake = image_patch.exists("cake")
        >>>     is_gummy_bear = image_patch.exists("gummy bear")
        >>>     return bool_to_yesno(is_cake and is_gummy_bear)
        """
        return len(self.find(object_name)) > 0

    def verify_property(self, object_name: str, property: str) -> bool:
        """Returns True if the object possesses the property, and False otherwise.
        Differs from 'exists' in that it presupposes the existence of the object specified by object_name, instead checking whether the object
        possesses the property.
        Parameters
```

```
116              -------
117              object_name : str
118                  A string describing the name of the object to be found in the image.
119              property : str
120                  A string describing the property to be checked.
121
122              Examples
123              -------
124              >>> # Do the letters have blue color?
125              >>> def execute_command(image) -> str:
126              >>>     image_patch = ImagePatch(image)
127              >>>     letters_patches = image_patch.find("letters")
128              >>>     # Question assumes only one letter patch
129              >>>     if len(letters_patches) == 0:
130              >>>         # If no letters are found, query the image directly
131              >>>         return image_patch.simple_query("Do the letters have blue color?")
132              >>>     return bool_to_yesno(letters_patches[0].verify_property("letters", "blue"))
133              """
134              return verify_property(self.cropped_image, object_name, property)
135
136          def best_text_match(self, option_list: List[str]) -> str:
137              """Returns the string that best matches the image.
138              Parameters
139              -------
140              option_list : str
141                  A list with the names of the different options
142              prefix : str
143                  A string with the prefixes to append to the options
144
145              Examples
146              -------
147              >>> # Is the cap gold or white?
148              >>> def execute_command(image)->str:
149              >>>     image_patch = ImagePatch(image)
150              >>>     cap_patches = image_patch.find("cap")
151              >>>     # Question assumes one cap patch
152              >>>     if len(cap_patches) == 0:
153              >>>         # If no cap is found, query the image directly
154              >>>         return image_patch.simple_query("Is the cap gold or white?")
155              >>>     return cap_patches[0].best_text_match(["gold", "white"])
156              """
157              return best_text_match(self.cropped_image, option_list)
158
159          def simple_query(self, question: str = None) -> str:
160              """Returns the answer to a basic question asked about the image. If no question is provided, returns the answer to "What is this?".
161              Parameters
162              -------
163              question : str
164                  A string describing the question to be asked.
165
166              Examples
167              -------
168
169              >>> # Which kind of animal is not eating?
170              >>> def execute_command(image) -> str:
171              >>>     image_patch = ImagePatch(image)
172              >>>     animal_patches = image_patch.find("animal")
173              >>>     for animal_patch in animal_patches:
174              >>>         if not animal_patch.verify_property("animal", "eating"):
175              >>>             return animal_patch.simple_query("What kind of animal is eating?") # crop would include eating so keep it in the query
176              >>>     # If no animal is not eating, query the image directly
177              >>>     return image_patch.simple_query("Which kind of animal is not eating?")
178
179              >>> # What is in front of the horse?
180              >>> # contains a relation (around, next to, on, near, on top of, in front of, behind, etc), so ask directly
181              >>> return image_patch.simple_query("What is in front of the horse?")
182              >>>
183              """
184              return simple_qa(self.cropped_image, question)
185
186          def compute_depth(self):
187              """Returns the median depth of the image crop
188              Parameters
189              ----------
190              Returns
191              -------
192              float
193                  the median depth of the image crop
194
```

```
195          Examples
196          --------
197          >>> # the person furthest away
198          >>> def execute_command(image)->ImagePatch:
199          >>>     image_patch = ImagePatch(image)
200          >>>     person_patches = image_patch.find("person")
201          >>>     person_patches.sort(key=lambda person: person.compute_depth())
202          >>>     return person_patches[-1]
203          """
204          depth_map = compute_depth(self.cropped_image)
205          return depth_map.median()
206
207      def crop(self, left: int, lower: int, right: int, upper: int) -> ImagePatch:
208          """Returns a new ImagePatch cropped from the current ImagePatch.
209          Parameters
210          -------
211          left : int
212              The leftmost pixel of the cropped image.
213          lower : int
214              The lowest pixel of the cropped image.
215          right : int
216              The rightmost pixel of the cropped image.
217          upper : int
218              The uppermost pixel of the cropped image.
219          -------
220          """
221          return ImagePatch(self.cropped_image, left, lower, right, upper)
222
223      def overlaps_with(self, left, lower, right, upper):
224          """Returns True if a crop with the given coordinates overlaps with this one,
225          else False.
226          Parameters
227          ----------
228          left : int
229              the left border of the crop to be checked
230          lower : int
231              the lower border of the crop to be checked
232          right : int
233              the right border of the crop to be checked
234          upper : int
235              the upper border of the crop to be checked
236
237          Returns
238          -------
239          bool
240              True if a crop with the given coordinates overlaps with this one, else False
241
242          Examples
243          --------
244          >>> # black cup on top of the table
245          >>> def execute_command(image) -> ImagePatch:
246          >>>     image_patch = ImagePatch(image)
247          >>>     table_patches = image_patch.find("table")
248          >>>     if len(table_patches) == 0:
249          >>>         table_patches = [image_patch]  # If no table found, assume the whole image is a table
250          >>>     table_patch = table_patches[0]
251          >>>     cup_patches = image_patch.find("black cup")
252          >>>     for cup in cup_patches:
253          >>>         if cup.vertical_center > table_patch.vertical_center
254          >>>             return cup
255          >>>     return cup_patches[0]  # If no cup found on top of the table, return the first cup found
256          """
257          return self.left <= right and self.right >= left and self.lower <= upper and self.upper >= lower
258
259
260  def best_image_match(list_patches: List[ImagePatch], content: List[str], return_index=False) -> Union[ImagePatch, int]:
261      """Returns the patch most likely to contain the content.
262      Parameters
263      ----------
264      list_patches : List[ImagePatch]
265      content : List[str]
266          the object of interest
267      return_index : bool
268          if True, returns the index of the patch most likely to contain the object
269
270      Returns
271      -------
272      int
273          Patch most likely to contain the object
```

```python
274
275         Examples
276         --------
277         >>> # Return the man with the hat
278         >>> def execute_command(image):
279         >>>     image_patch = ImagePatch(image)
280         >>>     man_patches = image_patch.find("man")
281         >>>     if len(man_patches) == 0:
282         >>>         return image_patch
283         >>>     hat_man = best_image_match(list_patches=man_patches, content=["hat"])
284         >>>     return hat_man
285
286         >>> # Return the woman with the pink scarf and blue pants
287         >>> def execute_command(image):
288         >>>     image_patch = ImagePatch(image)
289         >>>     woman_patches = image_patch.find("woman")
290         >>>     if len(woman_patches) == 0:
291         >>>         return image_patch
292         >>>     woman_most = best_image_match(list_patches=woman_patches, content=["pink scarf", "blue pants"])
293         >>>     return woman_most
294         """
295         return best_image_match(list_patches, content, return_index)


298 def distance(patch_a: ImagePatch, patch_b: ImagePatch) -> float:
299     """
300     Returns the distance between the edges of two ImagePatches. If the patches overlap, it returns a negative distance
301     corresponding to the negative intersection over union.
302     """
303     return distance(patch_a, patch_b)


306 def bool_to_yesno(bool_answer: bool) -> str:
307     return "yes" if bool_answer else "no"


310 def llm_query(question: str) -> str:
311     '''Answers a text question using GPT-3. The input question is always a formatted string with a variable in it.
312
313     Parameters
314     ----------
315     question: str
316         the text question to ask. Must not contain any reference to 'the image' or 'the photo', etc.
317     '''
318     return llm_query(question)


321 class VideoSegment:
322     """A Python class containing a set of frames represented as ImagePatch objects, as well as relevant information.
323     Attributes
324     ----------
325     video : torch.Tensor
326         A tensor of the original video.
327     start : int
328         An int describing the starting frame in this video segment with respect to the original video.
329     end : int
330         An int describing the ending frame in this video segment with respect to the original video.
331     num_frames->int
332         An int containing the number of frames in the video segment.
333
334     Methods
335     -------
336     frame_iterator->Iterator[ImagePatch]
337     trim(start, end)->VideoSegment
338         Returns a new VideoSegment containing a trimmed version of the original video at the [start, end] segment.
339     select_answer(info, question, options)->str
340         Returns the answer to the question given the options and additional information.
341     """
342
343     def __init__(self, video: torch.Tensor, start: int = None, end: int = None, parent_start=0, queues=None):
344         """Initializes a VideoSegment object by trimming the video at the given [start, end] times and stores the
345         start and end times as attributes. If no times are provided, the video is left unmodified, and the times are
346         set to the beginning and end of the video.
347
348         Parameters
349         -------
350         video : torch.Tensor
351             A tensor of the original video.
352         start : int
```

```
353              An int describing the starting frame in this video segment with respect to the original video.
354        end : int
355              An int describing the ending frame in this video segment with respect to the original video.
356        """
357
358        if start is None and end is None:
359            self.trimmed_video = video
360            self.start = 0
361            self.end = video.shape[0]  # duration
362        else:
363            self.trimmed_video = video[start:end]
364            if start is None:
365                start = 0
366            if end is None:
367                end = video.shape[0]
368            self.start = start + parent_start
369            self.end = end + parent_start
370
371        self.num_frames = self.trimmed_video.shape[0]
372
373    def frame_iterator(self) -> Iterator[ImagePatch]:
374        """Returns an iterator over the frames in the video segment."""
375        for i in range(self.num_frames):
376            yield ImagePatch(self.trimmed_video[i], self.start + i)
377
378    def trim(self, start: Union[int, None] = None, end: Union[int, None] = None) -> VideoSegment:
379        """Returns a new VideoSegment containing a trimmed version of the original video at the [start, end]
380        segment.
381
382        Parameters
383        ----------
384        start : Union[int, None]
385              An int describing the starting frame in this video segment with respect to the original video.
386        end : Union[int, None]
387              An int describing the ending frame in this video segment with respect to the original video.
388
389        Examples
390        --------
391        >>> # Return the second half of the video
392        >>> def execute_command(video):
393        >>>     video_segment = VideoSegment(video)
394        >>>     video_second_half = video_segment.trim(video_segment.num_frames // 2, video_segment.num_frames)
395        >>>     return video_second_half
396        """
397        if start is not None:
398            start = max(start, 0)
399        if end is not None:
400            end = min(end, self.num_frames)
401
402        return VideoSegment(self.trimmed_video, start, end, self.start)
403
404    def select_answer(self, info: dict, question: str, options: List[str]) -> str:
405        return select_answer(self.trimmed_video, info, question, options)
406
407    def __repr__(self):
408        return "VideoSegment({}, {})".format(self.start, self.end)
```
Listing 1. **Full API.**

Not all methods are used in all the benchmarks. Next we describe in more detail what content is used for the API specifications for every benchmark.

- **RefCOCO and RefCOCO+**. We use all the methods from the `ImagePatch` class except for `best_text_match` and `simple_query`. We also use the `best_text_match` and `distance` functions. Additionally we add `ImagePatch` usage examples in the API definition that are representative of the RefCOCO dataset, and look like the following:

```
1  # chair at the front
2  def execute_command(image) -> ImagePatch:
3      # Return the chair
4      image_patch = ImagePatch(image)
5      chair_patches = image_patch.find("chair")
6      chair_patches.sort(key=lambda chair: chair.compute_depth())
7      chair_patch = chair_patches[0]
8      # Remember: return the chair
9      return chair_patch
```
Listing 2. **RefCOCO example.**

- **GQA**. The GQA API contains all the contents in the API from Listing 1 up until the `llm_query` function, which is not used. The `ImagePatch` usage examples look like the following:

```python
# Is there a backpack to the right of the man?
def execute_command(image)->str:
    image_patch = ImagePatch(image)
    man_patches = image_patch.find("man")
    # Question assumes one man patch
    if len(man_patches) == 0:
        # If no man is found, query the image directly
        return image_patch.simple_query("Is there a backpack to the right of the man?")
    man_patch = man_patches[0]
    backpack_patches = image_patch.find("backpack")
    # Question assumes one backpack patch
    if len(backpack_patches) == 0:
        return "no"
    for backpack_patch in backpack_patches:
        if backpack_patch.horizontal_center > man_patch.horizontal_center:
            return "yes"
    return "no"
```
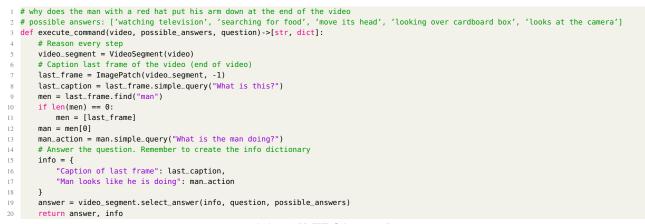
Listing 3. **GQA example.**

- **OK-VQA**. The API only uses the `simple_query` method from `ImagePatch`. It additionally uses the `llm_query` function. The `ImagePatch` usage examples look like the following:

```python

# Who is famous for allegedly doing this in a lightning storm?
def execute_command(image)->str:
    # The question is not direct perception, so we need to ask the image for more information
    # Salient information: what is being done?
    image = ImagePatch(image)
    guesses = []
    action = image.simple_query("What is being done?")
    external_knowledge_query = "Who is famous for allegedly {} in a lightning storm?".format(action)
    step_by_step_guess = llm_query(external_knowledge_query)
    guesses.append("what is being done is {}".format(action) + ", so " + step_by_step_guess)
    direct_guess = image.simple_query("Who is famous for allegedly doing this in a lightning storm?")
    guesses.append(direct_guess)
    return process_guesses("Who is famous for allegedly doing this in a lightning storm?", guesses)
```

Listing 4. **OK-VQA example.**

- **NeXT-QA**. The `VideoSegment` class is added to the API definition, and the available `ImagePatch` methods are `find`, `exists`, `best_text_match` and `simple_query`. The function `best_image_match` is also used. The `ImagePatch` usage examples look like:

```python
# why does the man with a red hat put his arm down at the end of the video
# possible answers: ['watching television', 'searching for food', 'move its head', 'looking over cardboard box', 'looks at the camera']
def execute_command(video, possible_answers, question)->[str, dict]:
    # Reason every step
    video_segment = VideoSegment(video)
    # Caption last frame of the video (end of video)
    last_frame = ImagePatch(video_segment, -1)
    last_caption = last_frame.simple_query("What is this?")
    men = last_frame.find("man")
    if len(men) == 0:
        men = [last_frame]
    man = men[0]
    man_action = man.simple_query("What is the man doing?")
    # Answer the question. Remember to create the info dictionary
    info = {
        "Caption of last frame": last_caption,
        "Man looks like he is doing": man_action
    }
    answer = video_segment.select_answer(info, question, possible_answers)
    return answer, info
```

Listing 5. **NeXT-QA example.**

- **Beyond benchmarks**. For the examples in Figure 1 we use the same API as the one used for the benchmarks, and the usage examples are taken from the benchmark APIs, combining them to have more generality. We do not add any other example, `ViperGPT` generalizes to the complex cases shown in Figure 1 just based on the provided API.

Note that in some of the examples we added comments, as well as error handling. The generated code also contains similar lines. We removed those for clarity in the figures shown in the main paper.