

# Qwen/Qwen2.5-VL-3B-Instruct

like 514

Follow Qwen 51.3k

Image-Text-to-Text Transformers Safetensors English qwen2\_5\_vl

image-to-text multimodal conversational text-generation-inference

arxiv:2309.00071 arxiv:2409.12191 arxiv:2308.12966

⋮

Train

Deploy

Use this model

Model card Files xet Community 41

Downloads last month  
4,090,190



Safetensors ⓘ

Model size 3.75B params Tensor type BF16 Chat template Files info

## Inference Providers NEW

Image-Text-to-Text

This model isn't deployed by any Inference Provider.

👤 7 Ask for provider support

## Model tree for Qwen/Qwen2.5-VL-3B-Instruct

Adapters	53 models
Finetunes	486 models
Quantizations	59 models

Spaces using Qwen/Qwen2.5-VL-3B-Instruct 100

## Collection including Qwen/Qwen2.5-VL-3B-Instruct

Qwen2.5-VL Collection

Vision-language model series based on Q... • 11 items • Updated Jul 21 • △ 538

Qwen2.5-VL-3B-Instruct

## 🔗 Introduction

In the past five months since Qwen2-VL's release, numerous developers have built new models on the Qwen2-VL vision-language models, providing us with valuable feedback. During this period, we focused on building more useful vision-language models. Today, we are excited to introduce the latest addition to the Qwen family: Qwen2.5-VL.

## 🔗 Key Enhancements:

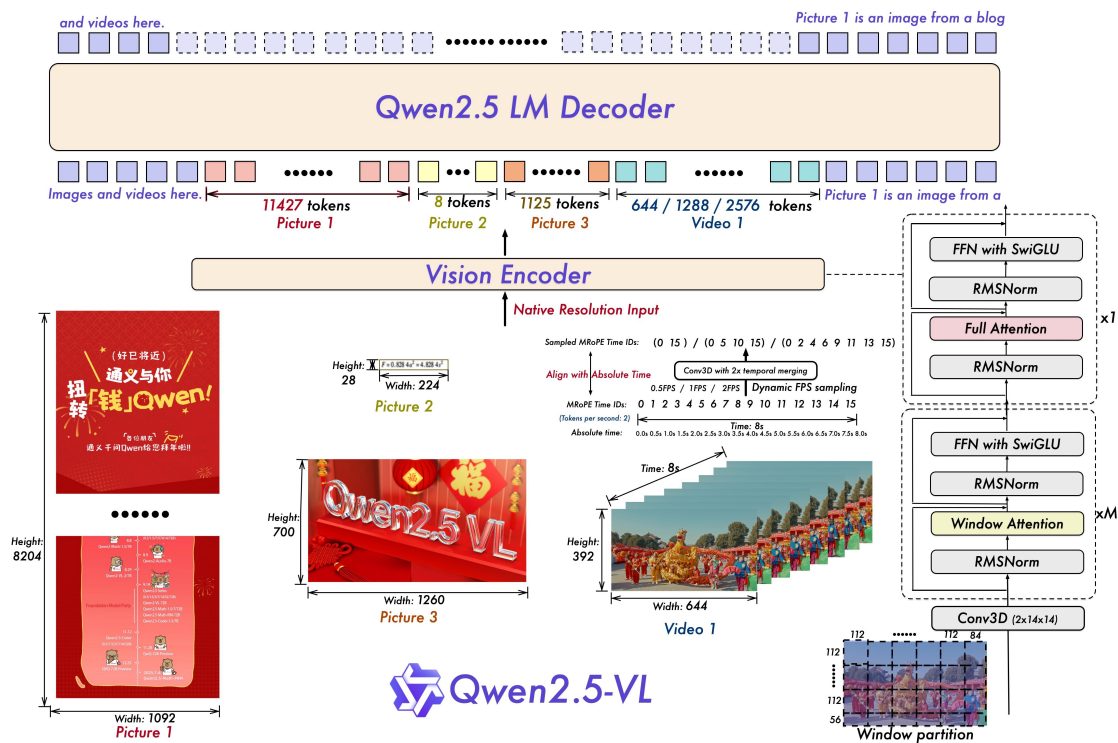
- **Understand things visually:** Qwen2.5-VL is not only proficient in recognizing common objects such as flowers, birds, fish, and insects, but it is highly capable of analyzing texts, charts, icons, graphics, and layouts within images.
- **Being agentic:** Qwen2.5-VL directly plays as a visual agent that can reason and dynamically direct tools, which is capable of computer use and phone use.
- **Understanding long videos and capturing events:** Qwen2.5-VL can comprehend videos of over 1 hour, and this time it has a new ability of capturing event by pinpointing the relevant video segments.
- **Capable of visual localization in different formats:** Qwen2.5-VL can accurately localize objects in an image by generating bounding boxes or points, and it can provide stable JSON outputs for coordinates and attributes.
- **Generating structured outputs:** for data like scans of invoices, forms, tables, etc. Qwen2.5-VL supports structured outputs of their contents, benefiting usages in finance, commerce, etc.

## 🔗 Model Architecture Updates:

- **Dynamic Resolution and Frame Rate Training for Video Understanding:**

We extend dynamic resolution to the temporal dimension by adopting dynamic FPS sampling, enabling the model to comprehend videos at various sampling rates. Accordingly, we update mRoPE

in the time dimension with IDs and absolute time alignment, enabling the model to learn temporal sequence and speed, and ultimately acquire the ability to pinpoint specific moments.



Streamlined and Efficient Vision Encoder

We enhance both training and inference speeds by strategically implementing window attention into the ViT. The ViT architecture is further optimized with SwiGLU and RMSNorm, aligning it with the structure of the Qwen2.5 LLM.

We have three models with 3, 7 and 72 billion parameters. This repo contains the instruction-tuned 3B Qwen2.5-VL model. For more information, visit our [Blog](#) and [GitHub](#).

Evaluation

Image benchmark

Benchmark	InternVL2.5-4B	Qwen2-VL-7B	Qwen2.5-VL-3B
MMMU <sub>val</sub>	52.3	54.1	53.1
MMMU-Pro <sub>val</sub>	32.7	30.5	31.6

Benchmark	InternVL2.5-4B	Qwen2-VL-7B	Qwen2.5-VL-3B
AI2D <sub>test</sub>	81.4	<b>83.0</b>	81.5
DocVQA <sub>test</sub>	91.6	94.5	<b>93.9</b>
InfoVQA <sub>test</sub>	72.1	76.5	<b>77.1</b>
TextVQA <sub>val</sub>	76.8	<b>84.3</b>	79.3
MMBench-V1.1 <sub>test</sub>	79.3	<b>80.7</b>	77.6
MMStar	58.3	<b>60.7</b>	55.9
MathVista <sub>testmini</sub>	60.5	58.2	<b>62.3</b>
MathVision <sub>full</sub>	20.9	16.3	<b>21.2</b>

[🔗](#) Video benchmark

Benchmark	InternVL2.5-4B	Qwen2-VL-7B	Qwen2.5-VL-3B
MVBench	71.6	67.0	67.0
VideoMME	63.6/62.3	69.0/63.3	67.6/61.5
MLVU	48.3	-	68.2
LVBench	-	-	43.3
MMBench-Video	1.73	1.44	1.63
EgoSchema	-	-	64.8
PerceptionTest	-	-	66.9
TempCompass	-	-	64.4
LongVideoBench	55.2	55.6	54.2
CharadesSTA/mIoU	-	-	38.8

## [🔗 Agent benchmark](#)

Benchmarks	Qwen2.5-VL-3B
ScreenSpot	55.5
ScreenSpot Pro	23.9
AITZ_EM	76.9
Android Control High_EM	63.7
Android Control Low_EM	22.2
AndroidWorld_SR	90.8
MobileMiniWob++_SR	67.9

## [🔗 Requirements](#)



The code of Qwen2.5-VL has been in the latest Hugging face transformers and we advise you to build from source with command:

```
pip install git+https://github.com/huggingface/transformers accelerate
```

or you might encounter the following error:

```
KeyError: 'qwen2_5_vl'
```

## [🔗 Quickstart](#)

Below, we provide simple examples to show how to use Qwen2.5-VL with  ModelScope and  Transformers.

The code of Qwen2.5-VL has been in the latest Hugging face transformers and we advise you to build from source with command:

```
pip install git+https://github.com/huggingface/transformers accelerate
```

or you might encounter the following error:

```
KeyError: 'qwen2_5_vl'
```

We offer a toolkit to help you handle various types of visual input more conveniently, as if you were using an API. This includes base64, URLs, and interleaved images and videos. You can install it using the following command:

```
# It's highly recommended to use '[decord]' feature for faster video loading.  
pip install qwen-vl-utils[decord]==0.0.8
```

If you are not using Linux, you might not be able to install `decord` from PyPI. In that case, you can use `pip install qwen-vl-utils` which will fall back to using torchvision for video processing. However, you can still [install decord from source](#) to get decord used when loading video.

## 🔗 Using 😊 Transformers to Chat

Here we show a code snippet to show you how to use the chat model with `transformers` and `qwen_vl_utils`:

```
from transformers import Qwen2_5_VLForConditionalGeneration, AutoTokenizer, AutoProcessor  
from qwen_vl_utils import process_vision_info  
  
# default: Load the model on the available device(s)  
model = Qwen2_5_VLForConditionalGeneration.from_pretrained(  
    "Qwen/Qwen2.5-VL-3B-Instruct", torch_dtype="auto", device_map="auto"  
)  
  
# We recommend enabling flash_attention_2 for better acceleration and memory saving
```

```

# model = Qwen2_5_VLForConditionalGeneration.from_pretrained(
#     "Qwen/Qwen2.5-VL-3B-Instruct",
#     torch_dtype=torch.bfloat16,
#     attn_implementation="flash_attention_2",
#     device_map="auto",
# )

# default processor
processor = AutoProcessor.from_pretrained("Qwen/Qwen2.5-VL-3B-Instruct")

# The default range for the number of visual tokens per image in the model is 4-16
# You can set min_pixels and max_pixels according to your needs, such as a token 1
# min_pixels = 256*28*28
# max_pixels = 1280*28*28
# processor = AutoProcessor.from_pretrained("Qwen/Qwen2.5-VL-3B-Instruct", min_pi

messages = [
    {
        "role": "user",
        "content": [
            {
                "type": "image",
                "image": "https://qianwen-res.oss-cn-beijing.aliyuncs.com/Qwen-VL/
            },
            {"type": "text", "text": "Describe this image."},
        ],
    }
]

# Preparation for inference
text = processor.apply_chat_template(
    messages, tokenize=False, add_generation_prompt=True
)
image_inputs, video_inputs = process_vision_info(messages)
inputs = processor(
    text=[text],
    images=image_inputs,
    videos=video_inputs,
    padding=True,
    return_tensors="pt",

```

```

)
inputs = inputs.to("cuda")

# Inference: Generation of the output
generated_ids = model.generate(**inputs, max_new_tokens=128)
generated_ids_trimmed = [
    out_ids[len(in_ids) :] for in_ids, out_ids in zip(inputs.input_ids, generated_ids)
]
output_text = processor.batch_decode(
    generated_ids_trimmed, skip_special_tokens=True, clean_up_tokenization_spaces=False
)
print(output_text)

```

- ▶ Multi image inference
- ▶ Video inference
- ▶ Batch inference

## [ModelScope](#)

We strongly advise users especially those in mainland China to use ModelScope.

`snapshot_download` can help you solve issues concerning downloading checkpoints.

## [More Usage Tips](#)

For input images, we support local files, base64, and URLs. For videos, we currently only support local files.

```

# You can directly insert a local file path, a URL, or a base64-encoded image into the prompt
## Local file path
messages = [
    {
        "role": "user",
        "content": [
            {"type": "image", "image": "file:///path/to/your/image.jpg"},
            {"type": "text", "text": "Describe this image."},
        ],
    }
]

```



```

### Image URL
messages = [
    {
        "role": "user",
        "content": [
            {"type": "image", "image": "http://path/to/your/image.jpg"},
            {"type": "text", "text": "Describe this image."},
        ],
    }
]

### Base64 encoded image
messages = [
    {
        "role": "user",
        "content": [
            {"type": "image", "image": "data:image;base64,/9j/..."},
            {"type": "text", "text": "Describe this image."},
        ],
    }
]

```

## [🔗](#) Image Resolution for performance boost

The model supports a wide range of resolution inputs. By default, it uses the native resolution for input, but higher resolutions can enhance performance at the cost of more computation. Users can set the minimum and maximum number of pixels to achieve an optimal configuration for their needs, such as a token count range of 256-1280, to balance speed and memory usage.

```

min_pixels = 256 * 28 * 28
max_pixels = 1280 * 28 * 28
processor = AutoProcessor.from_pretrained(
    "Qwen/Qwen2.5-VL-3B-Instruct", min_pixels=min_pixels, max_pixels=max_pixels
)

```

Besides, We provide two methods for fine-grained control over the image size input to the model:

1. Define `min_pixels` and `max_pixels`: Images will be resized to maintain their aspect ratio within the range of `min_pixels` and `max_pixels`.
2. Specify exact dimensions: Directly set `resized_height` and `resized_width`. These values will be rounded to the nearest multiple of 28.

```
# min_pixels and max_pixels
messages = [
    {
        "role": "user",
        "content": [
            {
                "type": "image",
                "image": "file:///path/to/your/image.jpg",
                "resized_height": 280,
                "resized_width": 420,
            },
            {"type": "text", "text": "Describe this image."},
        ],
    }
]

# resized_height and resized_width
messages = [
    {
        "role": "user",
        "content": [
            {
                "type": "image",
                "image": "file:///path/to/your/image.jpg",
                "min_pixels": 50176,
                "max_pixels": 50176,
            },
            {"type": "text", "text": "Describe this image."},
        ],
    }
]
```

The current `config.json` is set for context length up to 32,768 tokens. To handle extensive inputs exceeding 32,768 tokens, we utilize [YaRN](#), a technique for enhancing model length extrapolation, ensuring optimal performance on lengthy texts.

For supported frameworks, you could add the following to `config.json` to enable YaRN:

```
{
  ...,
  "type": "yarn",
  "mrope_section": [
    16,
    24,
    24
  ],
  "factor": 4,
  "original_max_position_embeddings": 32768
}
```

However, it should be noted that this method has a significant impact on the performance of temporal and spatial localization tasks, and is therefore not recommended for use.

At the same time, for long video inputs, since MRoPE itself is more economical with ids, the `max_position_embeddings` can be directly modified to a larger value, such as 64k.

## [🔗](#) Citation

If you find our work helpful, feel free to give us a cite.

```
@misc{qwen2.5-VL,
  title = {Qwen2.5-VL},
  url = {https://qwenlm.github.io/blog/qwen2.5-vl/},
  author = {Qwen Team},
  month = {January},
  year = {2025}
}
```

```
@article{Qwen2VL,
```

```
title={Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Ar  
author={Wang, Peng and Bai, Shuai and Tan, Sinan and Wang, Shijie and Fan, Zhiha  
journal={arXiv preprint arXiv:2409.12191},  
year={2024}  
}
```

```
@article{Qwen-VL,  
title={Qwen-VL: A Versatile Vision-Language Model for Understanding, Localizatio  
author={Bai, Jinze and Bai, Shuai and Yang, Shusheng and Wang, Shijie and Tan, S  
journal={arXiv preprint arXiv:2308.12966},  
year={2023}  
}
```