

# Qwen2.5-VL Technical Report

Qwen Team, Alibaba Group  
Feb 2025

Jaehyeon Son, Junhyun Kim

# Presenter



- Jaehyeon Son
- ML PhD
- Interests: Embodied AI and Robotics



- Junhyun Kim
- MS Robotics
- Interests : Robot Manipulation and Embodied AI

# Outline

- Timeline and Overview
- Model Architecture
- Pre-training and Post-training
- Experiments and Results
- Strengths
- Weaknesses and Limitations

# Timeline



- Sep 2024
- 2B / 8B / 72B
- MRoPE



- Mar 2025
- 7B
- Added audio modality
- TMRoPE
- Thinker-Talker



- Aug 2023
- Qwen-VL / Qwen-VL-Chat
- First attempt to integrate vision encoder on Qwen



- Feb 2025
- 3B / 7B / 32B / 72B
- Dynamic FPS sampling
- SwiGLU / RMSNorm
- Window Attention

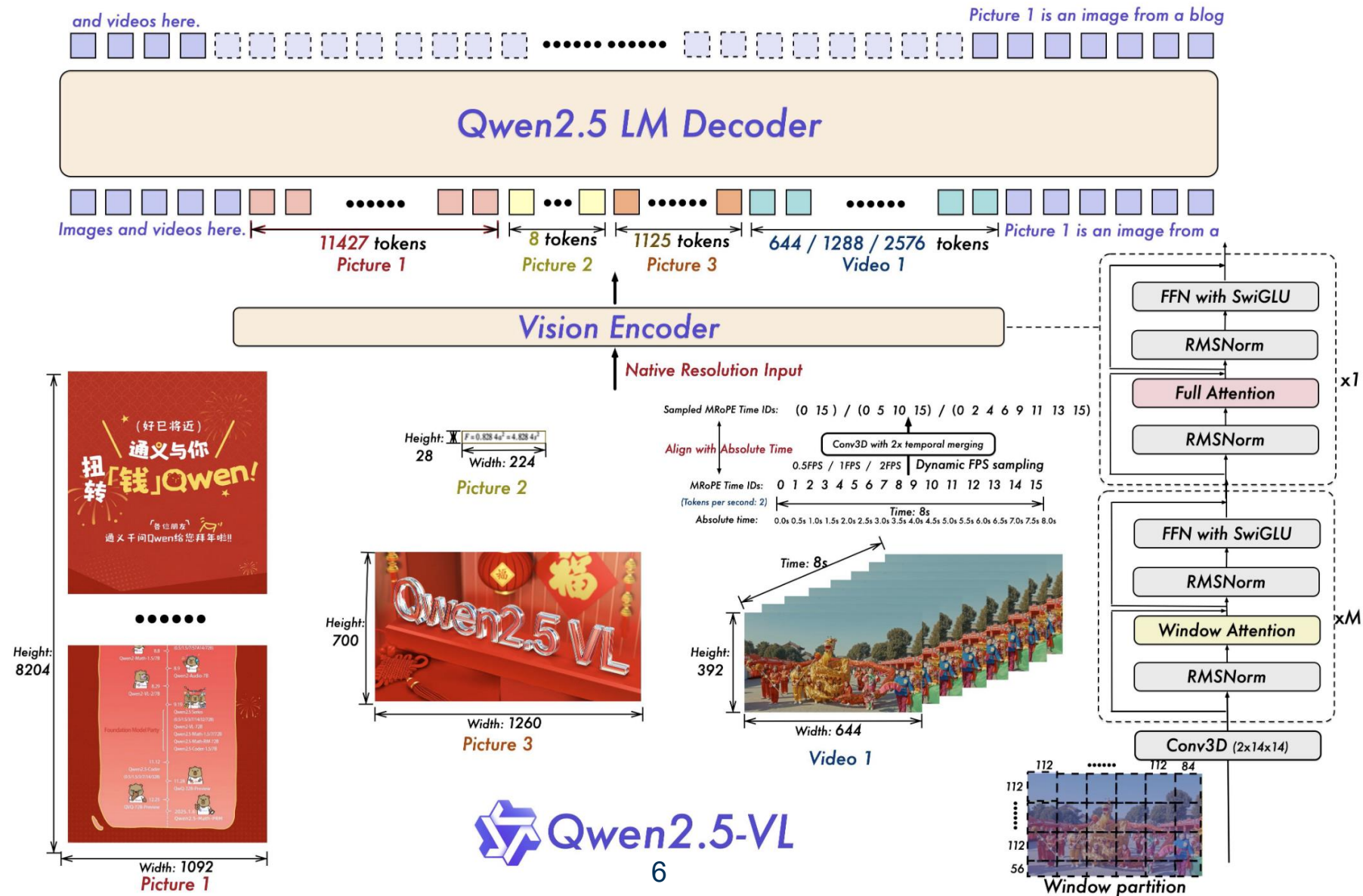


- **Sep 22, 2025**
- 30B
- Mixture-of-Experts(MoE)
- Audio Transformer(AuT)

# Overview Qwen2.5-VL

- Developer: Qwen Team, Alibaba
- Release: Feb 2025
- Model variants: 3B / 7B / 32B / 72B
- License: Apache 2.0 (✓ commercial use) / Qwen (✗ commercial use)
- Key features
  - Vision-language model for text, image, *video*
  - Long-video understanding (up to hours)
  - Multilingual support (English, Chinese, and many other languages)

# Model Architecture



# Image Encoding

- Image resizing: multiples of 28
- Support for dynamic resolution
  - Vanilla ViT<sup>[1]</sup>: fixed resolution only



Height: 28

Width: 224

$F = 0.8284a^2 = 4.8284r^2$

Picture 2

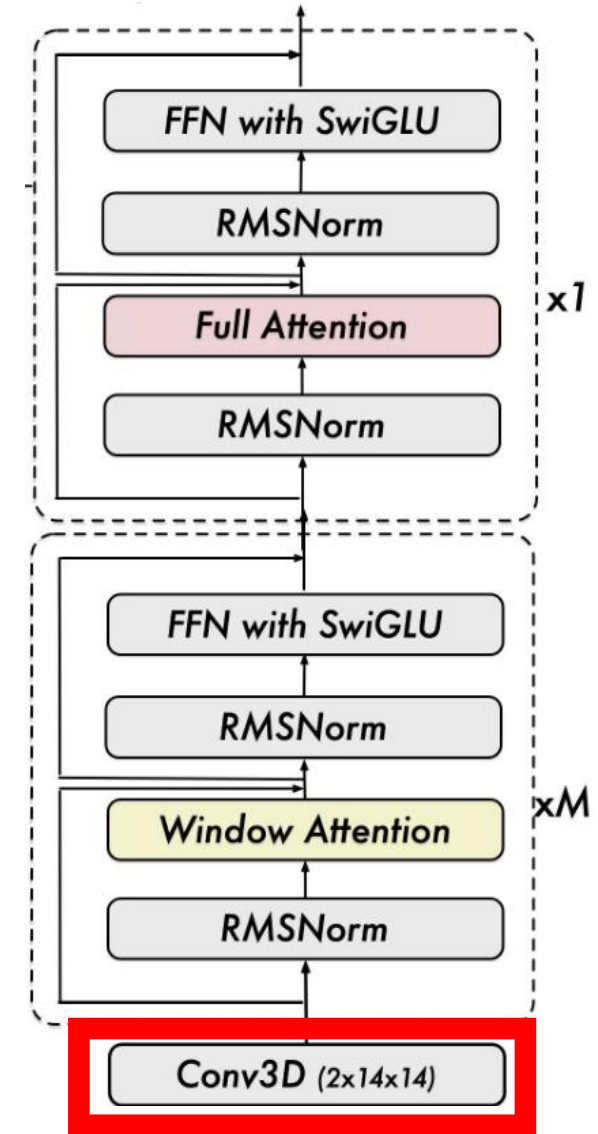


[1] Dosovitskiy et al., *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*, 2020



# Image Encoding

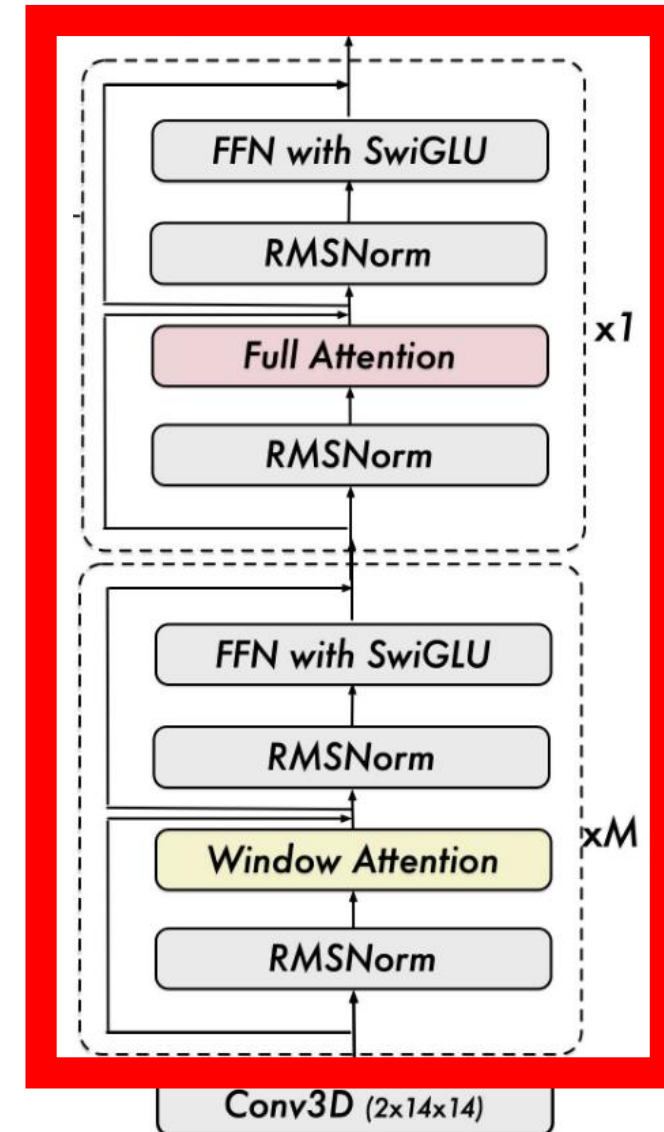
- Patch size: 14 x 14
  - Even number of patches in both width and height
- Each patch is passed to Conv3D





# Image Encoding

- Vision Encoder: ViT<sup>[1]</sup>
- Several modifications are applied
- After ViT: Vision-Language Merger
  - Merge 4 adjacent patch features into one token
  - Two-layer MLP



[1] Dosovitskiy et al., *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*, 2020

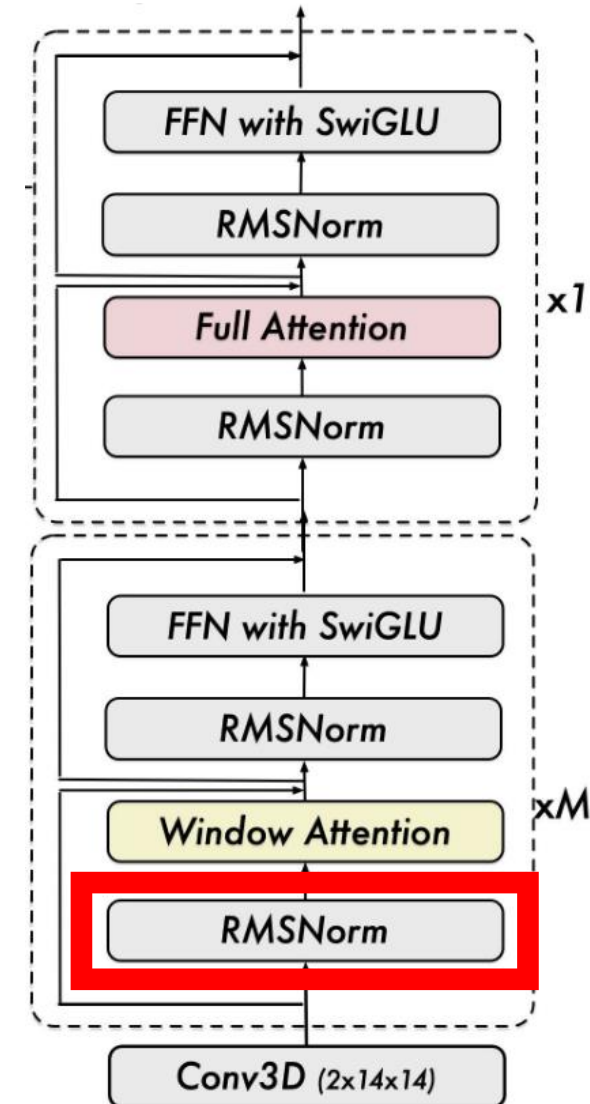
# Image Encoding

- RMSNorm<sup>[1]</sup>
  - Cheap & competitive alternative to vanilla LayerNorm
  - Employed in modern Transformers
    - T5, LLaMA, Mistral, etc.

$$\text{LayerNorm}(x) = \frac{x - \mu}{\sigma} \cdot \gamma + \beta$$

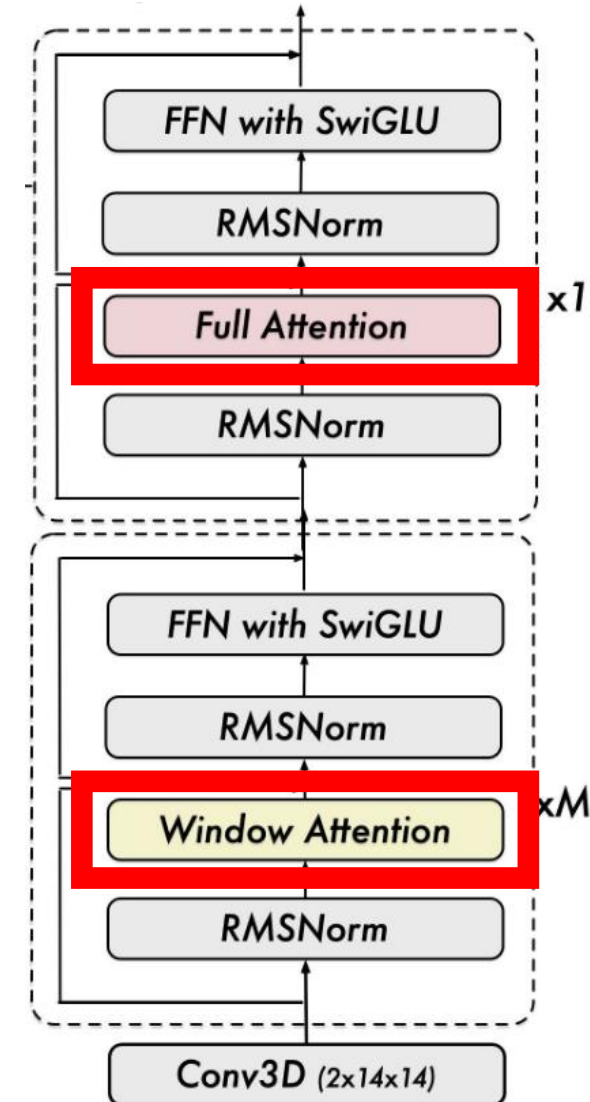
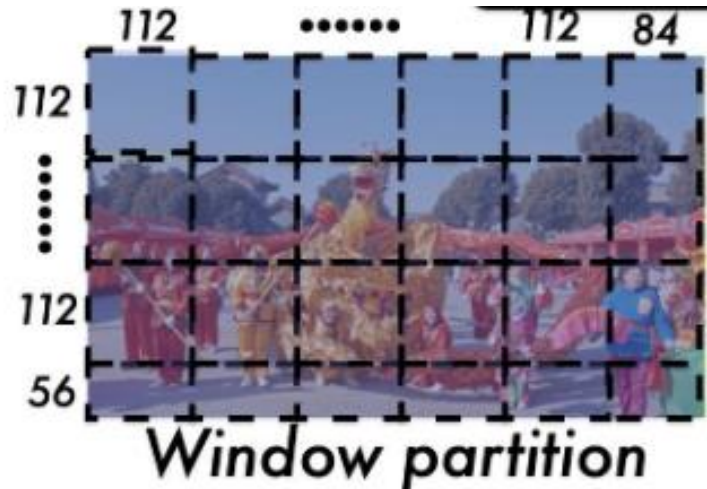
$$\text{RMSNorm}(x) = \frac{x}{\text{RMS}(x)} \cdot \gamma$$

$$\text{RMS}(x) = \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2 + \epsilon}$$



# Image Encoding

- Window Attention
  - Patch token attends only to adjacent 8×8 patches
    - 8x8 patches = 112x112 pixels
  - Linear complexity in the number of input tokens
- Full Attention is only used in 4 blocks



# Image Encoding

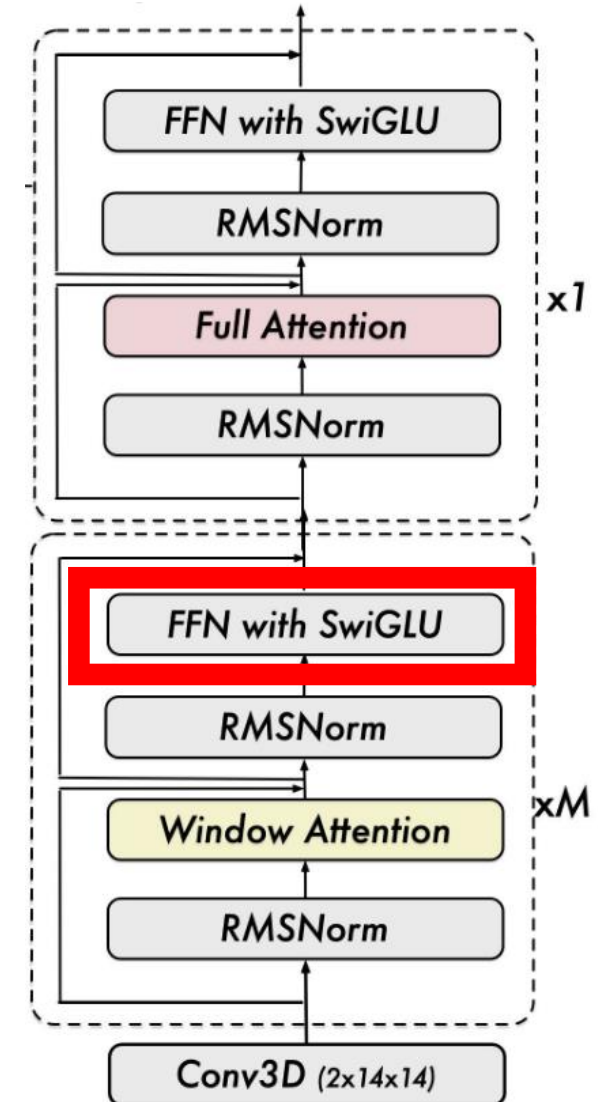
- FFN with SwiGLU<sup>[1]</sup>
  - Better expressivity than ReLU/GELU
  - Employed in modern Transformers
    - PaLM, LLaMA, Mistral, etc.

$$\text{FFN}(x) = W_2 f(W_1 x + b_1) + b_2$$

$$\text{GLU}(x) = (W_1 x) \odot \sigma(W_2 x)$$

$$\text{SwiGLU}(x) = (W_1 x) \odot \text{Swish}(W_2 x)$$

$$\text{Swish}(z) = z \cdot \sigma(z)$$



[1] Shazeer et al., *GLU Variants Improve Transformer*, 2020

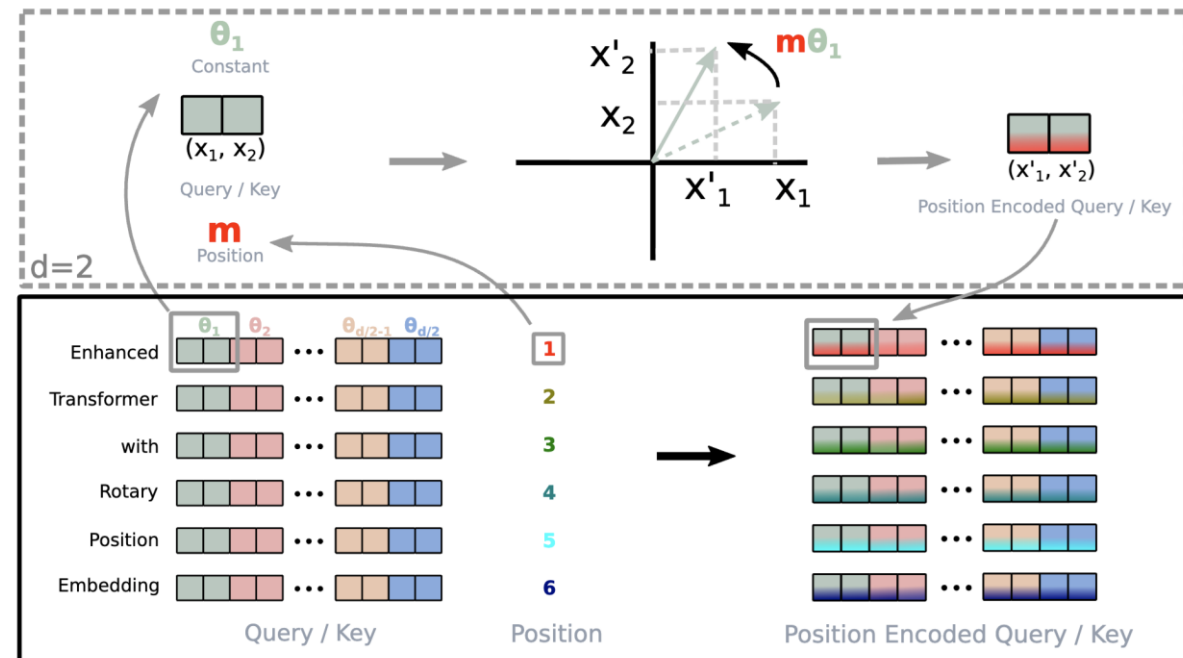
# Rotary Positional Embedding (RoPE<sup>[1]</sup>)

- Applies rotary transformations to Q/K
  - in *each* Attention layer
  - not to input embeddings
- Encodes relative distance directly in the dot product
- Preserves original token meaning through V

$$f_{\{q,k\}}(\mathbf{x}_m, m) = \mathbf{R}_{\Theta, m}^d \mathbf{W}_{\{q,k\}} \mathbf{x}_m$$

where

$$\mathbf{R}_{\Theta, m}^d = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix}$$



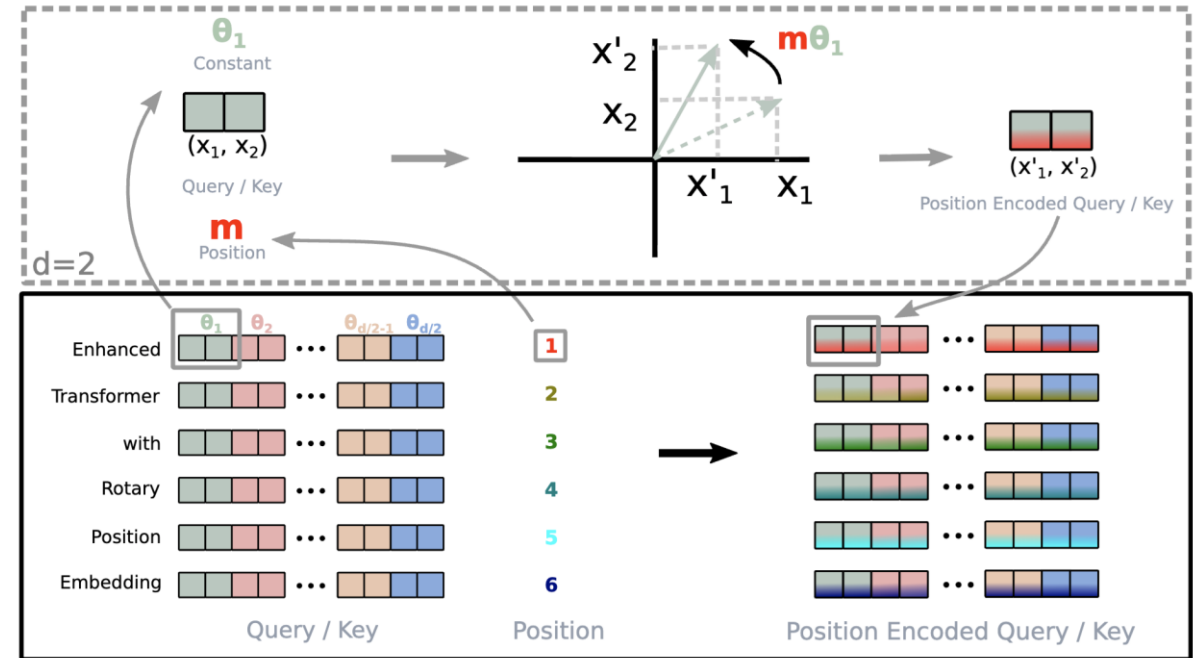
$$\Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, \dots, d/2]\}.$$

$$\mathbf{q}_m^\top \mathbf{k}_n = (\mathbf{R}_{\Theta, m}^d \mathbf{W}_q \mathbf{x}_m)^\top (\mathbf{R}_{\Theta, n}^d \mathbf{W}_k \mathbf{x}_n) = \mathbf{x}^\top \mathbf{W}_q \boxed{\mathbf{R}_{\Theta, n-m}^d} \mathbf{W}_k \mathbf{x}_n$$

# Rotary Positional Embedding (RoPE<sup>[1]</sup>)

- Image patch: 2D-RoPE
  - Halve the feature dimension
  - Encode row and column indices separately

Patch K/V



[1] Su et al., *RoFormer: Enhanced Transformer with Rotary Position Embedding*, 2021

# Rotary Positional Embedding (RoPE<sup>[1]</sup>)

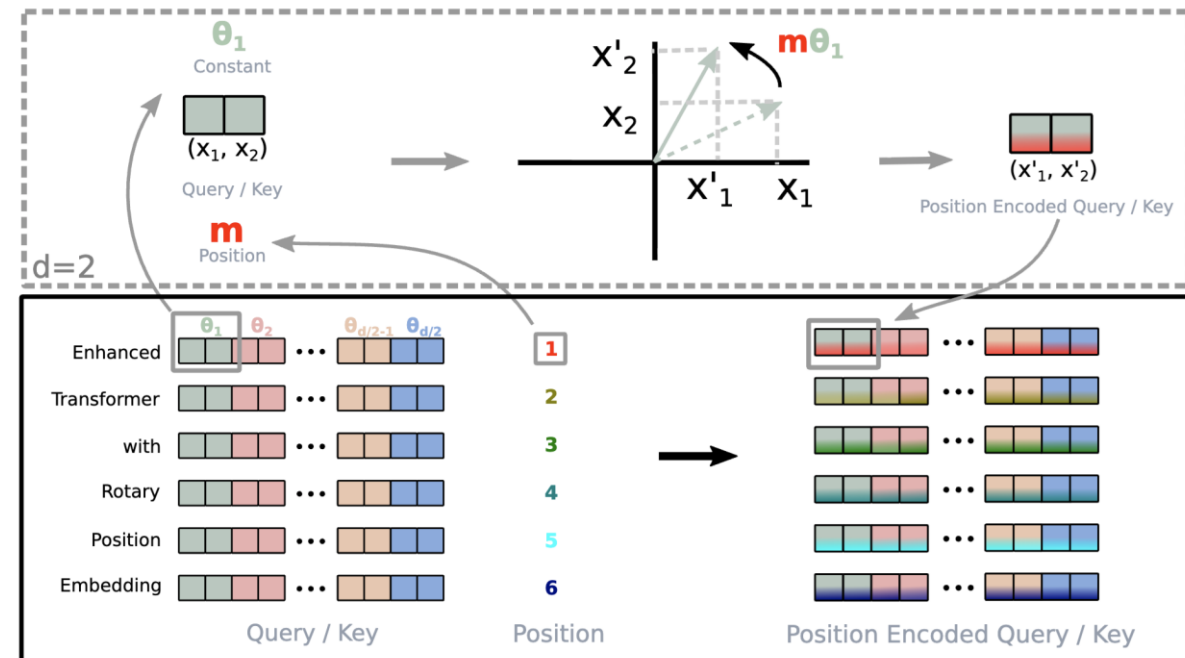
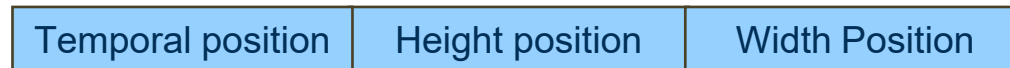
- Image patch: 2D-RoPE
  - Halve the feature dimension
  - Encode row and column indices separately

Patch K/V



- Video patch: Multimodal RoPE (MRoPE)

Patch K/V



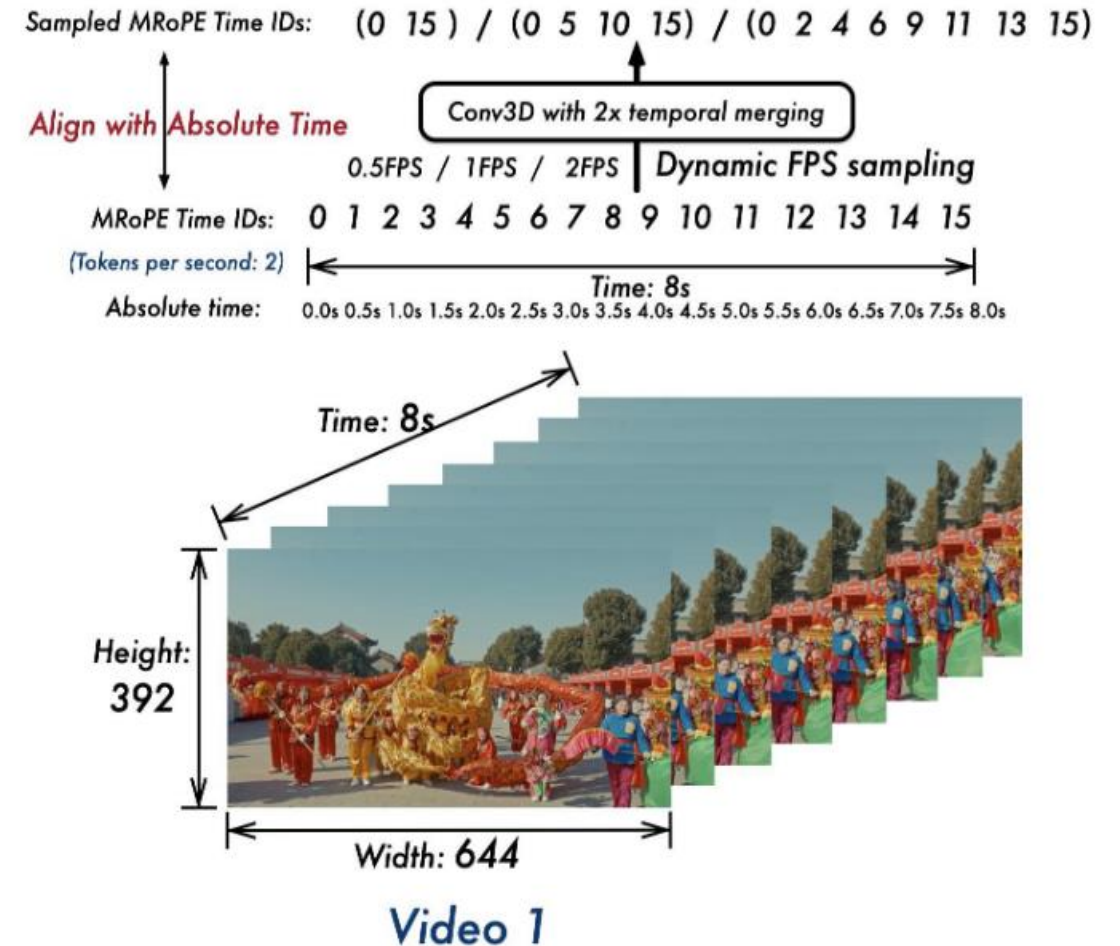
[1] Su et al., RoFormer: Enhanced Transformer with Rotary Position Embedding, 2021



# Video Encoding

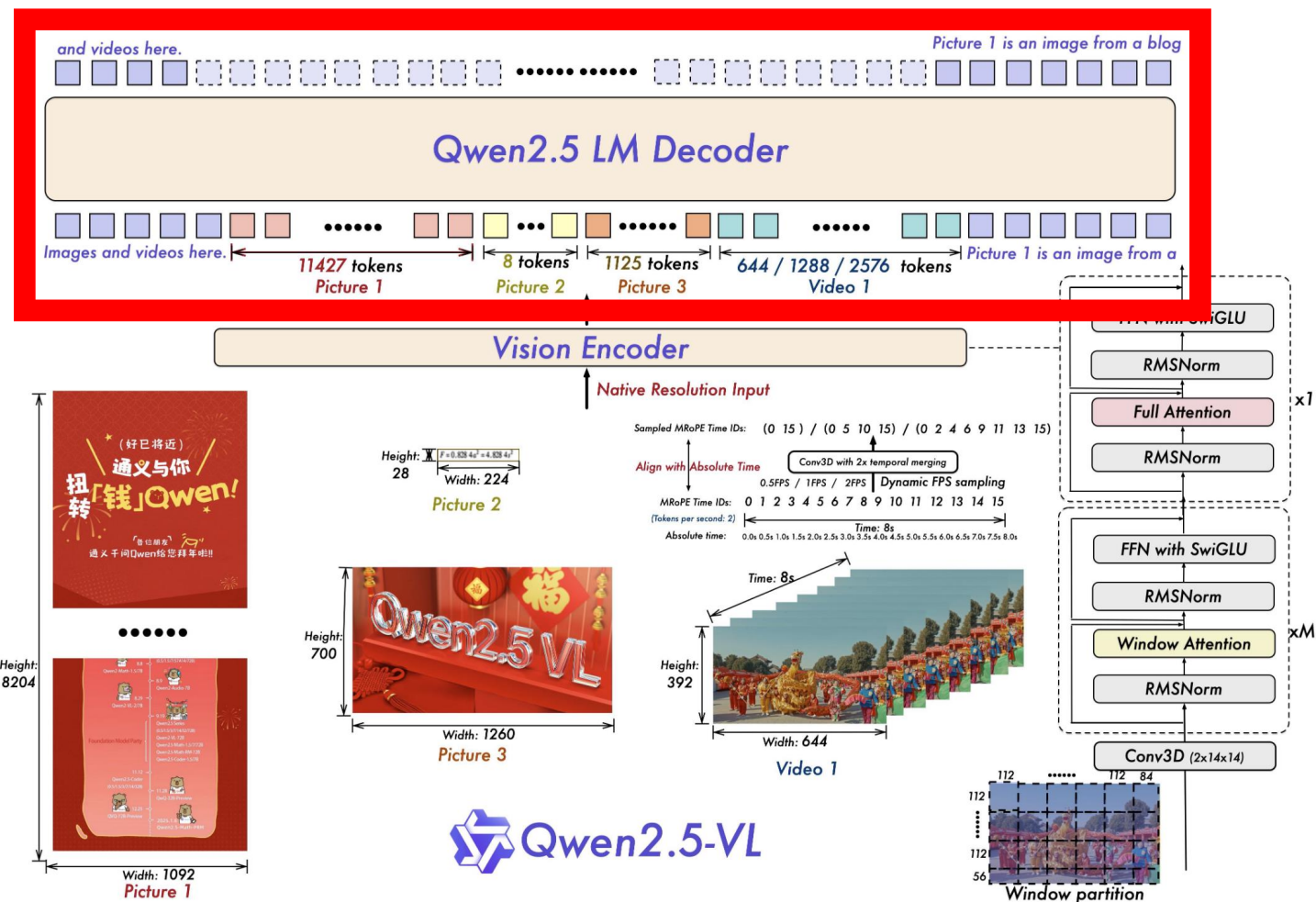
- Patch size: **2 x 14 x 14**
  - Bind two consecutive timesteps
- Dynamic FPS sampling
  - Vision Encoder handles various FPS
- Temporal position
  - aligned with absolute time
  - Qwen2-VL: simple frame index

Conv3D (2x14x14)



# LLM

- Initialized with Qwen2.5<sup>[1]</sup>
- We omit the details of LLM



[1] Qwen Team, Qwen2.5 Technical Report, 2025

# Model Configuration

Configuration	Qwen2.5-VL-3B	Qwen2.5-VL-7B	Qwen2.5-VL-72B
<b>Vision Transformer (ViT)</b>			
Hidden Size	1280	1280	1280
# Layers	32	32	32
# Num Heads	16	16	16
Intermediate Size	3456	3456	3456
Patch Size	14	14	14
Window Size	112	112	112
Full Attention Block Indexes	{7, 15, 23, 31}	{7, 15, 23, 31}	{7, 15, 23, 31}
<b>Vision-Language Merger</b>			
In Channel	1280	1280	1280
Out Channel	2048	3584	8192
<b>Large Language Model (LLM)</b>			
Hidden Size	2048	3,584	8192
# Layers	36	28	80
# KV Heads	2	4	8
Head Size	128	128	128
Intermediate Size	4864	18944	29568
Embedding Tying	✓	✗	✗
Vocabulary Size	151646	151646	151646
# Trained Tokens	4.1T	4.1T	4.1T

# Training Overview

## Pre-Training

1. Visual Pre-Training
2. Multimodal Pre-Training
3. Long-Context Pre-Training

## Post-Training

1. Supervised Fine-Tuning
  - Domain-Tailored Filtering
  - Rejection Sampling
2. Direct Preference Optimization

# Data Processing Overview

## Pre-Training

Interleaved  
Image-Text Data

Document  
Omni-Parsing Data

OCR Data

Video Data

Agent Data

## Post-Training

Instruction Data

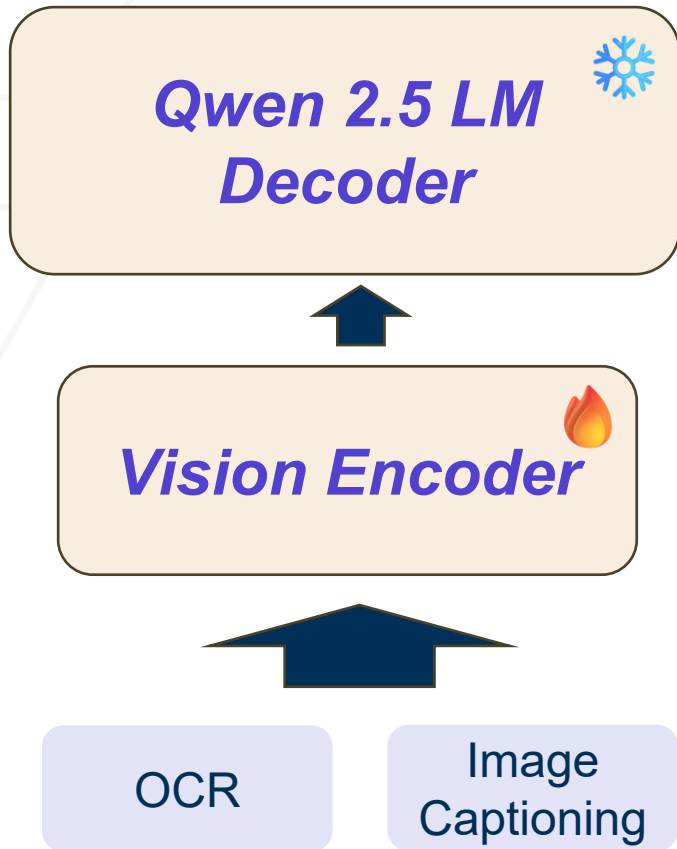
VQA  
Data

Image  
Captioning

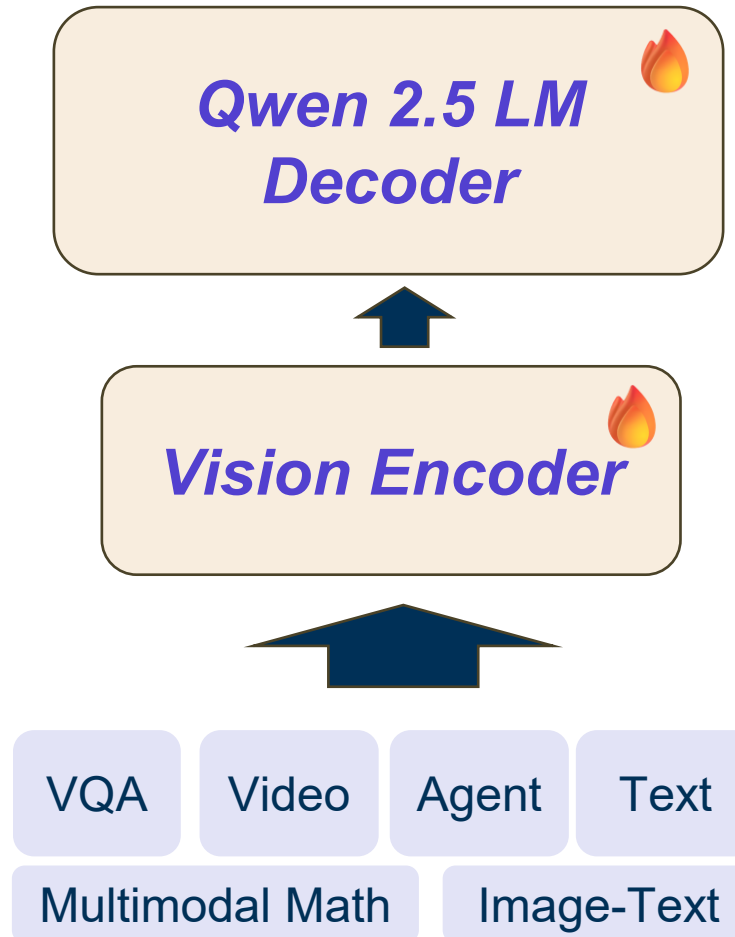
...

# Pre-Training

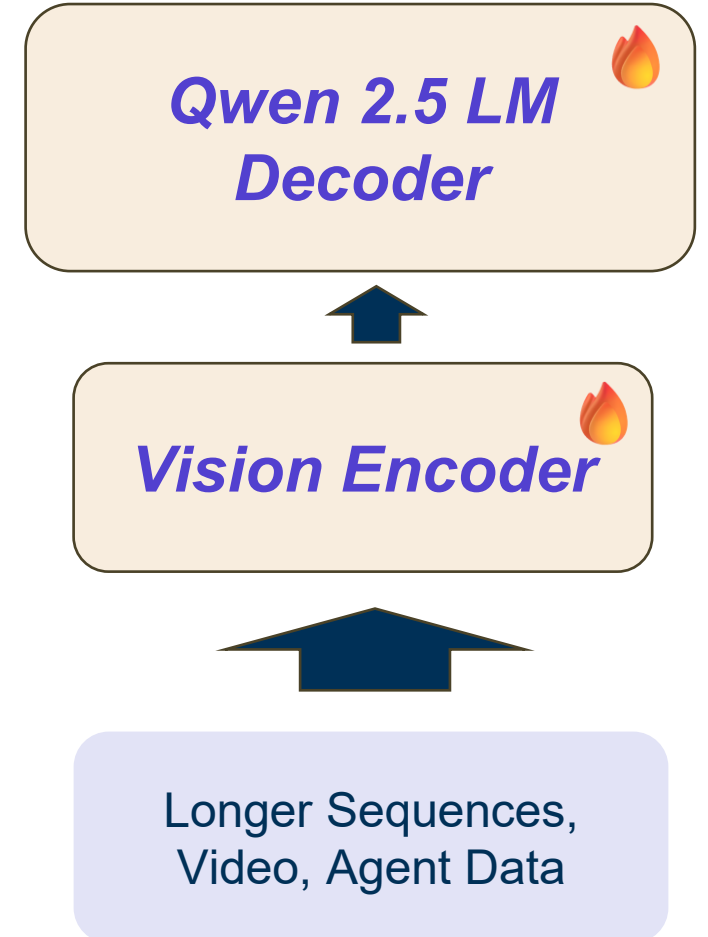
## 1. Visual Pre-Training



## 2. Multimodal Pre-Training



## 3. Long-Context Pre-Training



## Interleaved Image-Text Data

- Standard data cleaning (Wei Li et al. 2024)
- Four-stage scoring system
  - (1) text-only quality
  - (2) image-text relevance
  - (3) image-text complementarity
  - (4) information density balance



# Pre-Training Data

Interleaved  
Image-Text

Document  
Omni-Parsing

OCR

Video

Agent

## Document Omni-Parsing

- Convert document data into HTML format
- Include the coordinates for each module

### QwenVL HTML Format

```
<html><body>
# paragraph
<p data-bbox="x1 y1 x2 y2"> content </p>
# table
<style>table{id} style</style><table data-bbox="x1 y1 x2 y2" class="table{id}"> table content
</table>
# chart
<div class="chart" data-bbox="x1 y1 x2 y2"> <img data-bbox="x1 y1 x2 y2" /><table> chart content
</table></div>
# formula
<div class="formula" data-bbox="x1 y1 x2 y2"> <img data-bbox="x1 y1 x2 y2" /> <div> formula
content </div></div>
# image caption
<div class="image caption" data-bbox="x1 y1 x2 y2"> <img data-bbox="x1 y1 x2 y2" /><p> image
caption </p></div>
# image ocr
<div class="image ocr" data-bbox="x1 y1 x2 y2"> <img data-bbox="x1 y1 x2 y2" /><p> image ocr
</p></div>
# music sheet
<div class="music sheet" format="abc notation" data-bbox="x1 y1 x2 y2"> <img data-bbox="x1 y1
x2 y2" /> <div> music sheet content </div></div>
# chemical formula content
<div class="chemical formula" format="smile" data-bbox="x1 y1 x2 y2"> <img data-bbox="x1 y1
x2 y2" /> <div> chemical formula content </div></div>
</html></body>
```

# Pre-Training Data

Interleaved  
Image-Text

Document  
Omni-Parsing

OCR

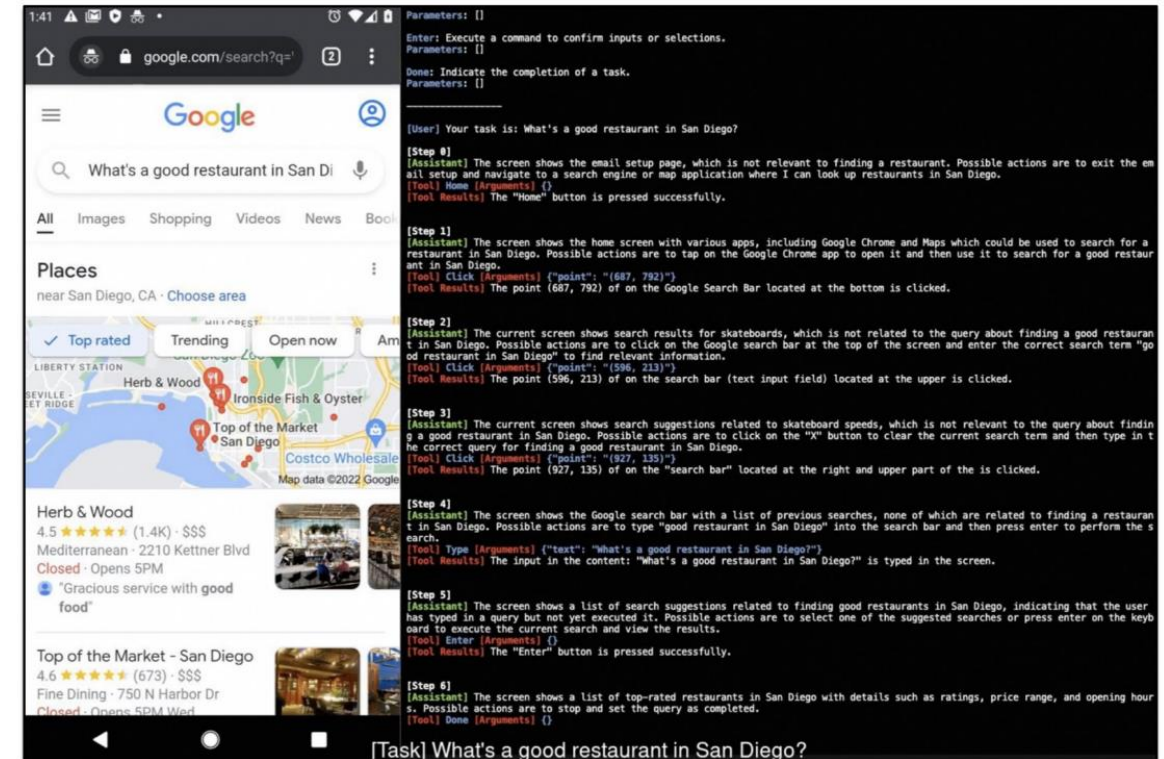
Video

Agent

## Agent Data

- Collect screenshots on mobile, web, and desktop platforms
- For decision-making, unify the operations across platforms into a function call
- Given a GT operation, explanation for reasoning content is generated by human and model annotators

## UI Interaction



# Pre-Training

Stages	Visual Pre-Training	Multimodal Pre-Training	Long-Context Pre-Training
Data	Image Caption Knowledge OCR	+ Pure text Interleaved Data VQA, Video Grounding, Agent	+ Long Video Long Agent Long Document
Tokens	1.5T	2T	0.6T
Sequence length	8192	8192	32768
Training	ViT	ViT & LLM	ViT & LLM

# Post-Training Data

## 1. *Supervised Fine-Tuning(SFT)*

- 2 million entries
- Evenly distributed between text data(50%) and multimodal data(50%)

VQA

Image  
Captioning

Math Problem  
Solving

Coding  
Tasks

Document

OCR

Video

Agent

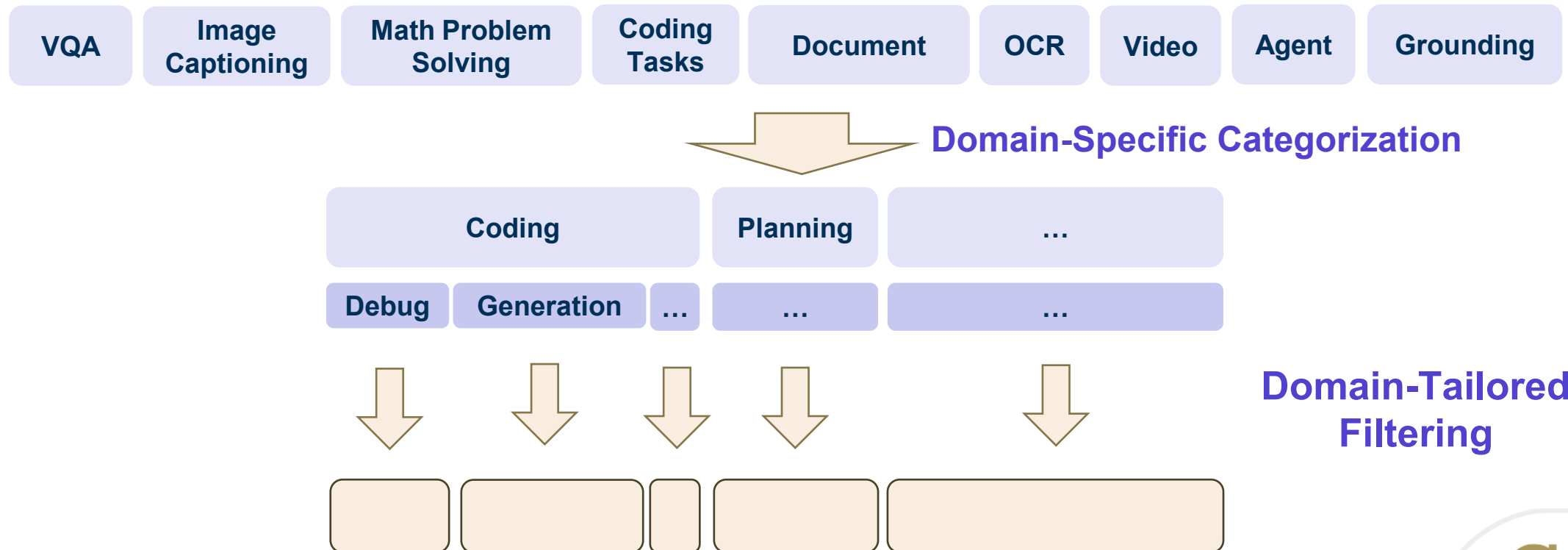
Grounding

## 2. Direct Preference Optimization(DPO)

# Post-Training Data

## 1. Supervised Fine-Tuning(SFT)

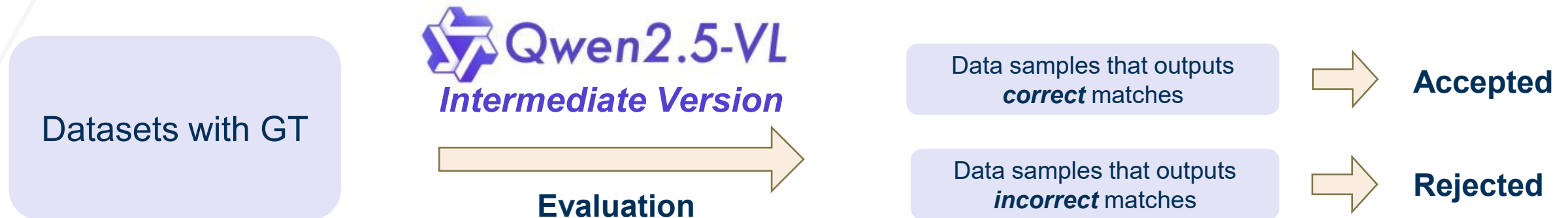
- 2 million entries
- Evenly distributed between text data(50%) and multimodal data(50%)



# Post-Training Data

## 1. Supervised Fine-Tuning(SFT)

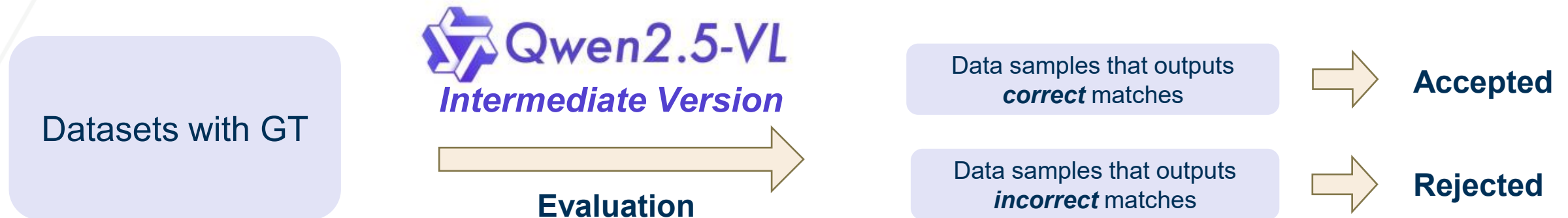
- Domain-Specific Categorization and Domain-Tailored Filtering
- ***Rejection Sampling***



# Post-Training Data

## 1. Supervised Fine-Tuning(SFT)

- Domain-Specific Categorization and Domain-Tailored Filtering
- ***Rejection Sampling***



**Does it learn easy dialogues  
and problems only?**



# Post-Training Data

1. Supervised Fine-Tuning(SFT)
  - Domain-Specific Categorization and Domain-Tailored Filtering
  - Rejection Sampling
  
2. ***Direct Preference Optimization(DPO)***
  - Only on image-text and pure text data

# Experiments and Results

Table 3: Performance of Qwen2.5-VL and State-of-the-art.

Datasets	Previous Open-source SoTA	Claude-3.5 Sonnet-0620	GPT-4o 0513	InternVL2.5 78B	Qwen2-VL 72B	Qwen2.5-VL 72B	Qwen2.5-VL 7B	Qwen2.5-VL 3B
<i>College-level Problems</i>								
MMMU <sub>val</sub> (Yue et al., 2023)	70.1 Chen et al. (2024d)	68.3	69.1	70.1	64.5	<b>70.2</b>	58.6	53.1
MMMU-Pro <sub>overall</sub> (Yue et al., 2024)	48.6 Chen et al. (2024d)	51.5	<b>51.9</b>	48.6	46.2	51.1	38.3	31.56
<i>Math</i>								
MathVista <sub>mini</sub> (Lu et al., 2024)	72.3 Chen et al. (2024d)	67.7	63.8	72.3	70.5	<b>74.8</b>	68.2	62.3
MATH-Vision <sub>full</sub> (Wang et al., 2024d)	32.2 Chen et al. (2024d)	-	30.4	32.2	25.9	<b>38.1</b>	25.1	21.2
MathVerse <sub>mini</sub> (Zhang et al., 2024c)	51.7 Chen et al. (2024d)	-	50.2	51.7	-	<b>57.6</b>	49.2	47.6
<i>General Visual Question Answering</i>								
MegaBench (Chen et al., 2024b)	47.4 MiniMax et al. (2025)	52.1	<b>54.2</b>	45.6	46.8	51.3	36.8	28.9
MMBench-EN <sub>test</sub> (Liu et al., 2023d)	88.3 Chen et al. (2024d)	82.6	83.4	88.3	86.9	<b>88.6</b>	83.5	79.1
MMBench-CN <sub>test</sub> (Liu et al., 2023d)	88.5 Chen et al. (2024d)	83.5	82.1	<b>88.5</b>	86.7	87.9	83.4	78.1
MMBench-V1.1-EN <sub>test</sub> (Liu et al., 2023d)	87.4 Chen et al. (2024d)	80.9	83.1	87.4	86.1	<b>88.4</b>	82.6	77.4
MMStar (Chen et al., 2024c)	69.5 Chen et al. (2024d)	65.1	64.7	69.5	68.3	<b>70.8</b>	63.9	55.9
MME <sub>sum</sub> (Fu et al., 2023)	<b>2494</b> Chen et al. (2024d)	1920	2328	<b>2494</b>	2483	2448	2347	2157
MuirBench (Wang et al., 2024a)	63.5 Chen et al. (2024d)	-	68.0	63.5	-	<b>70.7</b>	59.6	47.7
BLINK <sub>val</sub> (Fu et al., 2024c)	63.8 Chen et al. (2024d)	-	<b>68.0</b>	63.8	-	64.4	56.4	47.6
CRPE <sub>relation</sub> (Wang et al., 2024h)	78.8 Chen et al. (2024d)	-	76.6	78.8	-	<b>79.2</b>	76.4	73.6
HallBench <sub>avg</sub> (Guan et al., 2023)	<b>58.1</b> Wang et al. (2024f)	55.5	55.0	57.4	<b>58.1</b>	55.2	52.9	46.3
MTVQA (Tang et al., 2024)	<b>31.9</b> Chen et al. (2024d)	25.7	27.8	31.9	30.9	31.7	29.2	24.8
RealWorldQA <sub>avg</sub> (X.AI, 2024)	78.7 Chen et al. (2024d)	60.1	75.4	<b>78.7</b>	77.8	75.7	68.5	65.4
MME-RealWorld <sub>en</sub> (Zhang et al., 2024f)	62.9 Chen et al. (2024d)	51.6	45.2	62.9	-	<b>63.2</b>	57.4	53.1
MMVet <sub>turbo</sub> (Yu et al., 2024)	74.0 Wang et al. (2024f)	70.1	69.1	72.3	74.0	<b>76.2</b>	67.1	61.8
MM-MT-Bench (Agrawal et al., 2024)	7.4 Agrawal et al. (2024)	7.5	<b>7.72</b>	-	6.59	7.6	6.3	5.7

# Experiments and Results

Table 3: Performance of Qwen2.5-VL and State-of-the-art.

Datasets	Previous Open-source SoTA	Claude-3.5 Sonnet-0620	GPT-4o 0513	InternVL2.5 78B	Qwen2-VL 72B	Qwen2.5-VL 72B	Qwen2.5-VL 7B	Qwen2.5-VL 3B
<i>College-level Problems</i>						Performance degradation		
MMMU <sub>val</sub> (Yue et al., 2023)	70.1 Chen et al. (2024d)	68.3	69.1	70.1	64.5	70.2	58.6	53.1
MMMU-Pro <sub>overall</sub> (Yue et al., 2024)	48.6 Chen et al. (2024d)	51.5	51.9	48.6	46.2	51.1	38.3	31.56
<i>Math</i>								
MathVista <sub>mini</sub> (Lu et al., 2024)	72.3 Chen et al. (2024d)	67.7	63.8	72.3	70.5	74.8	68.2	62.3
MATH-Vision <sub>full</sub> (Wang et al., 2024d)	32.2 Chen et al. (2024d)	-	30.4	32.2	25.9	38.1	25.1	21.2
MathVerse <sub>mini</sub> (Zhang et al., 2024c)	51.7 Chen et al. (2024d)	-	50.2	51.7	-	57.6	49.2	47.6
<i>General Visual Question Answering</i>								
MegaBench (Chen et al., 2024b)	47.4 MiniMax et al. (2025)	52.1	54.2	45.6	46.8	51.3	36.8	28.9
MMBench-EN <sub>test</sub> (Liu et al., 2023d)	88.3 Chen et al. (2024d)	82.6	83.4	88.3	86.9	88.6	83.5	79.1
MMBench-CN <sub>test</sub> (Liu et al., 2023d)	88.5 Chen et al. (2024d)	83.5	82.1	88.5	86.7	87.9	83.4	78.1
MMBench-V1.1-EN <sub>test</sub> (Liu et al., 2023d)	87.4 Chen et al. (2024d)	80.9	83.1	87.4	86.1	88.4	82.6	77.4
MMStar (Chen et al., 2024c)	69.5 Chen et al. (2024d)	65.1	64.7	69.5	68.3	70.8	63.9	55.9
MME <sub>sum</sub> (Fu et al., 2023)	2494 Chen et al. (2024d)	1920	2328	2494	2483	2448	2347	2157
MuirBench (Wang et al., 2024a)	63.5 Chen et al. (2024d)	-	68.0	63.5	-	70.7	59.6	47.7
BLINK <sub>val</sub> (Fu et al., 2024c)	63.8 Chen et al. (2024d)	-	68.0	63.8	-	64.4	56.4	47.6
CRPE <sub>relation</sub> (Wang et al., 2024h)	78.8 Chen et al. (2024d)	-	76.6	78.8	-	79.2	76.4	73.6
HallBench <sub>avg</sub> (Guan et al., 2023)	58.1 Wang et al. (2024f)	55.5	55.0	57.4	58.1	55.2	52.9	46.3
MTVQA (Tang et al., 2024)	31.9 Chen et al. (2024d)	25.7	27.8	31.9	30.9	31.7	29.2	24.8
RealWorldQA <sub>avg</sub> (X.AI, 2024)	78.7 Chen et al. (2024d)	60.1	75.4	78.7	77.8	75.7	68.5	65.4
MME-RealWorld <sub>en</sub> (Zhang et al., 2024f)	62.9 Chen et al. (2024d)	51.6	45.2	62.9	-	63.2	57.4	53.1
MMVet <sub>turbo</sub> (Yu et al., 2024)	74.0 Wang et al. (2024f)	70.1	69.1	72.3	74.0	76.2	67.1	61.8
MM-MT-Bench (Agrawal et al., 2024)	7.4 Agrawal et al. (2024)	7.5	7.72	-	6.59	7.6	6.3	5.7

# Experiments and Results

Table 4: Performance on pure text tasks of the 70B+ Instruct models and Qwen2.5-VL.

Datasets	Llama-3.1-70B	Llama-3.1-405B	Qwen2-72B	Qwen2.5-72B	Qwen2.5-VL-72B
<i>General Tasks</i>					
MMLU-Pro	66.4	<b>73.3</b>	64.4	71.1	71.2
MMLU-redux	83.0	86.2	81.6	<b>86.8</b>	85.9
LiveBench-0831	46.6	53.2	41.5	52.3	<b>57.0</b>
<i>Mathematics &amp; Science Tasks</i>					
GPQA	46.7	<b>51.1</b>	42.4	49.0	49.0
MATH	68.0	73.8	69.0	<b>83.1</b>	83.0
GSM8K	95.1	<b>96.8</b>	93.2	95.8	95.3
<i>Coding Tasks</i>					
HumanEval	80.5	<b>89.0</b>	86.0	86.6	87.8
MultiPL-E	68.2	73.5	69.2	75.1	<b>79.5</b>
<i>Alignment Tasks</i>					
IFEval	83.6	86.0	77.6	84.1	<b>86.3</b>

# Experiments and Results

Retained performance on text-only evaluation, due to pure text data(50%) in train datasets

Table 4: Performance on pure text tasks of the 70B+ Inst

Datasets	Llama-3.1-70B	Llama-3.1-405B	Qwen2-72B	Qwen2.5-72B	Qwen2.5-VL-72B
<i>General Tasks</i>					
MMLU-Pro	66.4	<b>73.3</b>	64.4	71.1	71.2
MMLU-redux	83.0	86.2	81.6	<b>86.8</b>	85.9
LiveBench-0831	46.6	53.2	41.5	52.3	<b>57.0</b>
<i>Mathematics &amp; Science Tasks</i>					
GPQA	46.7	<b>51.1</b>	42.4	49.0	49.0
MATH	68.0	73.8	69.0	<b>83.1</b>	83.0
GSM8K	95.1	<b>96.8</b>	93.2	95.8	95.3
<i>Coding Tasks</i>					
HumanEval	80.5	<b>89.0</b>	86.0	86.6	87.8
MultiPL-E	68.2	73.5	69.2	75.1	<b>79.5</b>
<i>Alignment Tasks</i>					
IFEval	83.6	86.0	77.6	84.1	<b>86.3</b>

# Experiments and Results

Table 5: Performance of Qwen2.5-VL and other models on OCR, chart, and document understanding benchmarks.

Datasets	Claude-3.5 Sonnet	Gemini 1.5 Pro	GPT 4o	InternVL2.5 78B	Qwen2.5-VL 72B	Qwen2.5-VL 7B	Qwen2.5-VL 3B
<i>OCR-related Parsing Tasks</i>							
CC-OCR	62.5	73.0	66.9	64.7	<b>79.8</b>	77.8	74.5
OmniDocBench <sub>edit en/zh</sub> ↓	0.330/0.381	0.230/ <b>0.281</b>	0.265/0.435	0.275/0.324	<b>0.226/0.324</b>	0.308/0.398	0.409/0.543
<i>OCR-related Understanding Tasks</i>							
AI2D <sub>w. M.</sub>	81.2	88.4	84.6	<b>89.1</b>	88.7	83.9	81.6
TextVQA <sub>val</sub>	76.5	78.8	77.4	83.4	83.5	<b>84.9</b>	79.3
DocVQA <sub>test</sub>	95.2	93.1	91.1	95.1	<b>96.4</b>	95.7	93.9
InfoVQA <sub>test</sub>	74.3	81.0	80.7	84.1	<b>87.3</b>	82.6	77.1
ChartQA <sub>test Avg.</sub>	<b>90.8</b>	87.2	86.7	88.3	89.5	87.3	84.0
CharXiv <sub>RQ/DQ</sub>	<b>60.2/84.3</b>	43.3/72.0	47.1/84.5	42.4/82.3	49.7/ <b>87.4</b>	42.5/73.9	31.3/58.6
SEED-Bench-2-Plus	71.7	70.8	72.0	71.3	<b>73.0</b>	70.4	67.6
OCRBench	788	754	736	854	<b>885</b>	864	797
VCR <sub>En-Hard-EM</sub>	41.7	28.1	73.2	-	79.8	<b>80.5</b>	37.5
<i>OCR-related Comprehensive Tasks</i>							
OCRBench_v2 <sub>en/zh</sub>	45.2/39.6	51.9/43.1	46.5/32.2	49.8/52.1	<b>61.5/63.7</b>	56.3/57.2	54.3/52.1



# Experiments and Results

Table 6: Performance of Qwen2.5-VL and other models on grounding.

Datasets	Gemini 1.5 Pro	Grounding DINO	Molmo 72B	InternVL2.5 78B	Qwen2.5-VL 72B	Qwen2.5-VL 7B	Qwen2.5-VL 3B
Refcoco <sub>val</sub>	73.2	90.6	-	93.7	92.7	90.0	89.1
Refcoco <sub>testA</sub>	72.9	93.2	-	95.6	94.6	92.5	91.7
Refcoco <sub>testB</sub>	74.6	88.2	-	92.5	89.7	85.4	84.0
Refcoco+ <sub>val</sub>	62.5	88.2	-	90.4	88.9	84.2	82.4
Refcoco+ <sub>testA</sub>	63.9	89.0	-	94.7	92.2	89.1	88.0
Refcoco+ <sub>testB</sub>	65.0	75.9	-	86.9	83.7	76.9	74.1
Refcocog <sub>val</sub>	75.2	86.1	-	92.7	89.9	87.2	85.2
Refcocog <sub>test</sub>	76.2	87.0	-	92.2	90.3	87.2	85.7
ODinW	36.7	55.0	-	31.7	43.1	37.3	37.5
PointGrounding	-	-	69.2	-	67.5	67.3	58.3

Table 7: Performance of Qwen2.5-VL and other models on counting.

Datasets	Gemini 1.5-Pro	GPT-4o	Claude-3.5 Sonnet	Molmo-72b	InternVL2.5-78B	Qwen2.5-VL-72B
CountBench	85.5	87.9	89.7	91.2	72.1	93.6



# Experiments and Results

Significantly outperforming results on long video understanding, due to MRoPE and dynamic FPS sampling

Table 8: Performance of Qwen2.5-VL and other models on video benchmarks.

Long  
Video

Datasets	Gemini 1.5-Pro	GPT-4o	Qwen2.5-VL-72B	Qwen2.5-VL-7B	Qwen2.5-VL-3B
<i>Video Understanding Tasks</i>					
Video-MME <sub>w/o sub.</sub>	<b>75.0</b>	71.9	73.3	65.1	61.5
Video-MME <sub>w sub.</sub>	<b>81.3</b>	77.2	79.1	71.6	67.6
Video-MMMU	53.9	<b>61.2</b>	60.2	47.4	-
MMVU <sub>val</sub>	65.4	<b>67.4</b>	62.9	50.1	-
MVBench	60.5	64.6	<b>70.4</b>	69.6	67.0
MMBench-Video	1.30	1.63	<b>2.02</b>	1.79	1.63
LongVideoBench <sub>val</sub>	64.0	<b>66.7</b>	60.7	56.0	54.2
LVBench	33.1	30.8	<b>47.3</b>	45.3	43.3
EgoSchema <sub>test</sub>	71.2	72.2	<b>76.2</b>	65.0	64.8
PerceptionTest <sub>test</sub>	-	-	<b>73.2</b>	70.5	66.9
MLVU <sub>M-Avg</sub>	-	64.6	<b>74.6</b>	70.2	68.2
TempCompass <sub>Avg</sub>	67.1	73.8	<b>74.8</b>	71.7	64.4
<i>Video Grounding Tasks</i>					
Charades-STA <sub>mIoU</sub>	-	35.7	<b>50.9</b>	43.6	38.8

# Experiments and Results

Table 9: Performance of Qwen2.5-VL and other models on GUI Agent benchmarks.

Benchmarks	GPT-4o	Gemini 2.0	Claude	Aguvis-72B	Qwen2-VL-72B	Qwen2.5-VL-72B
ScreenSpot	18.1	84.0	83.0	<b>89.2</b>	-	87.1
ScreenSpot Pro	-	-	17.1	23.6	1.6	<b>43.6</b>
Android Control High <sub>EM</sub>	20.8	28.5	12.5	66.4	59.1	<b>67.36</b>
Android Control Low <sub>EM</sub>	19.4	60.2	19.4	84.4	59.2	<b>93.7</b>
AndroidWorld <sub>SR</sub>	34.5% (SoM)	26% (SoM)	27.9%	26.1%	6% (SoM)	<b>35%</b>
MobileMiniWob++ <sub>SR</sub>	61%	42% (SoM)	61% (SoM)	66%	50% (SoM)	<b>68%</b>
OSWorld	5.03	4.70	<b>14.90</b>	10.26	2.42	8.83

Outperform other models even when  
Set-of-Mark(SoM) is applied to their inputs

# Strengths

- Can handle images at their native resolution
- Improved long video comprehension by dynamic FPS sampling
- SOTA on some document parsing benchmark
- Strong visual agent capabilities
- Improved computational efficiency by applying RMSNorm and SwiGLU

# Weaknesses & Limitations

- Lacks long-form video performance comparison with the previous version, Qwen 2-VL

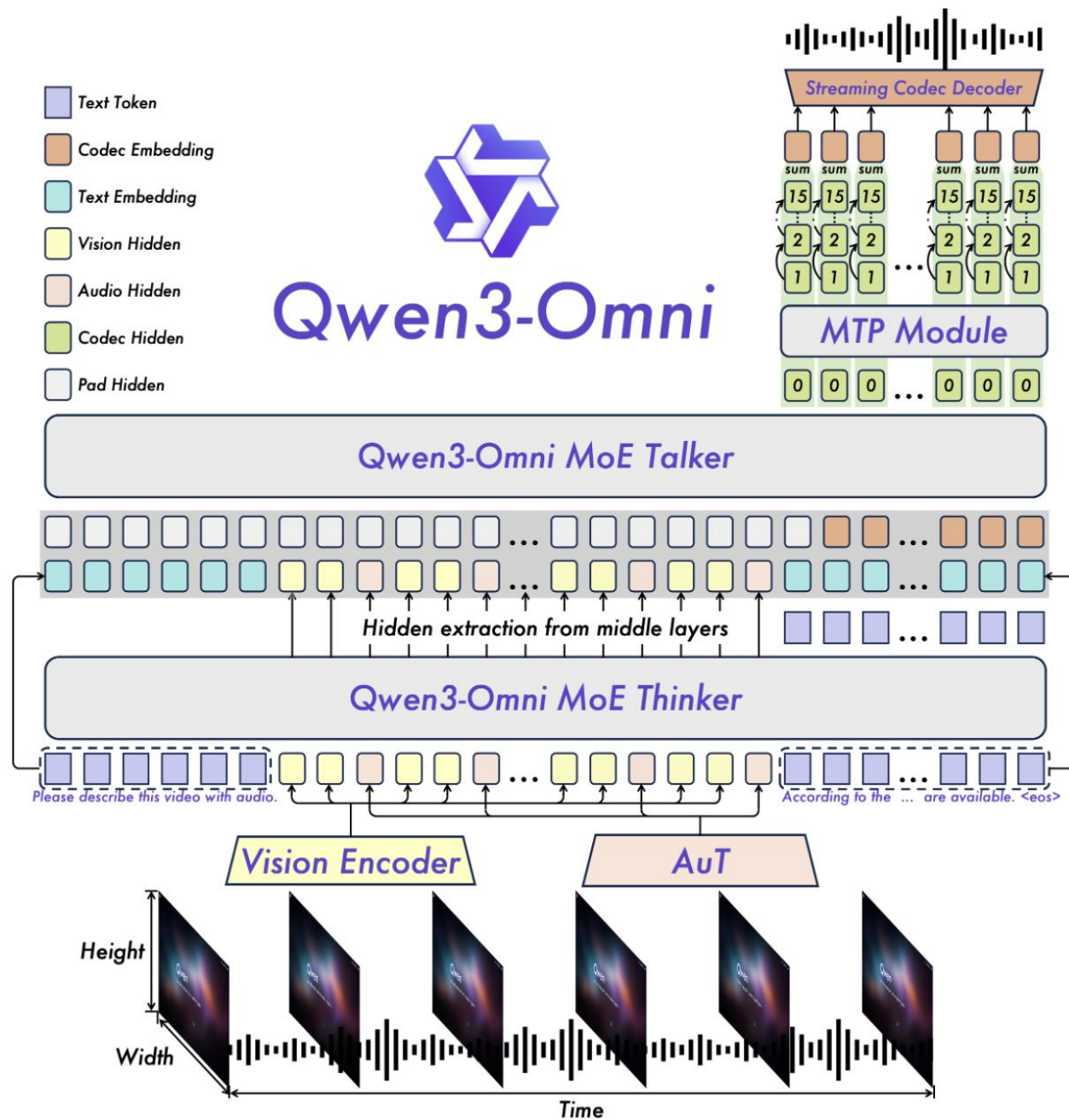
# Qwen3-Omni Technical Report

Qwen Team, Alibaba Group  
Sep 2025

Jaehyeon Son, Junhyun Kim

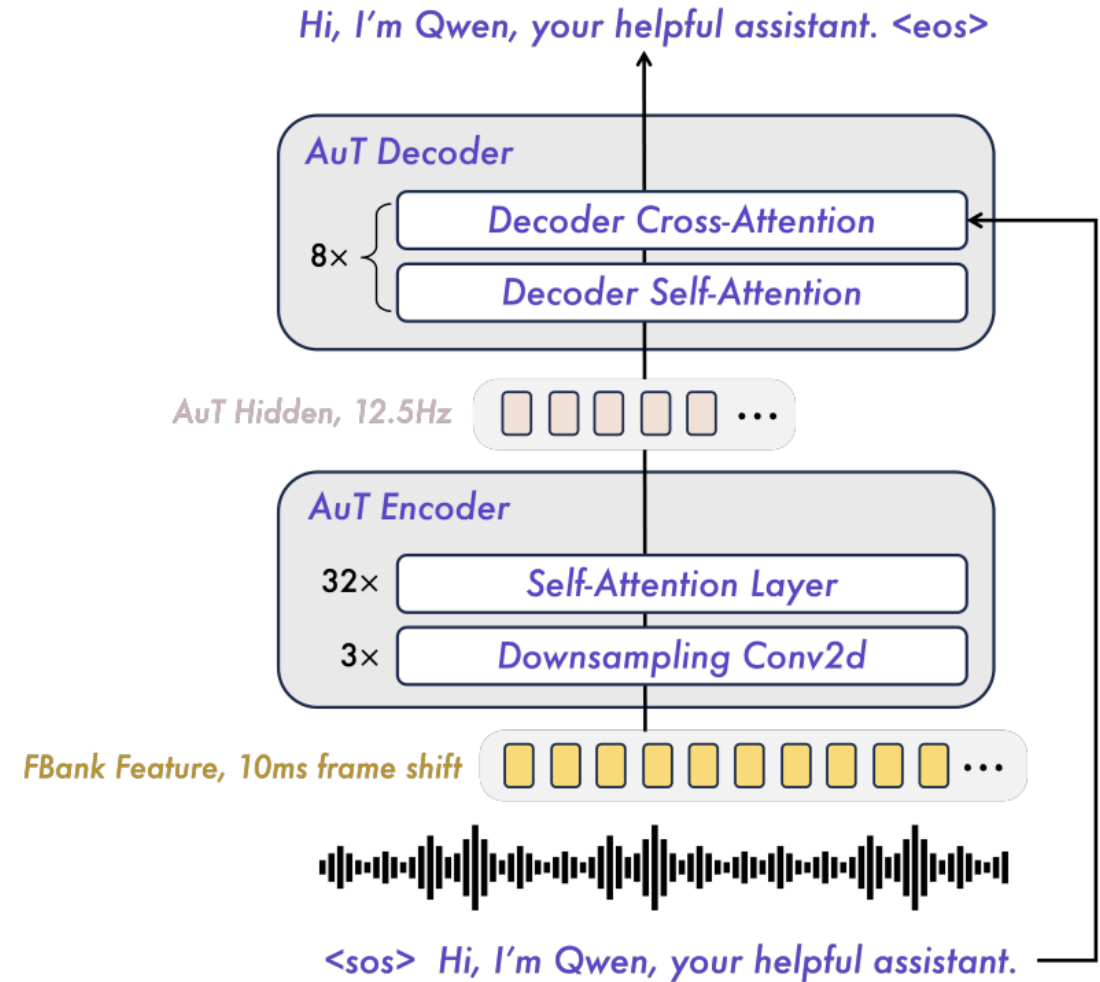
# Architecture

- Audio Transformer (AuT)
- Mixture of Experts (MoE)
- Thinker-Talker (Qwen 2.5-Omni)
- TMRoPE (Time-aligned MRoPE) (Qwen 2.5-Omni)



# Architecture

- Audio Transformer (AuT)
  - The filter bank features are downsampled using Conv2D blocks.
  - Token rate is reduced to 12.5 Hz in hidden layer.
  - 0.6B parameters



# Results

Table 11: AudioVisual  $\rightarrow$  Text performance of Qwen3-Omni-Instruct and other non-reasoning baselines. The highest scores are shown in bold.

Datasets	Previous Open-source SoTA	Gemini-2.5-Flash	Qwen2.5-Omni	Qwen3-Omni-30B-A3B -Instruct	Qwen3-Omni-Flash -Instruct
WorldSense	47.1(Yang et al., 2025b)	50.9	45.4	54.0	<b>54.1</b>

Table 12: AudioVisual  $\rightarrow$  Text performance of Qwen3-Omni-30B-A3B-Thinking and other reasoning baselines. The highest scores are shown in bold.

Datasets	Previous Open-source SoTA	Gemini-2.5-Flash -Thinking	Qwen3-Omni-30B-A3B -Thinking	Qwen3-Omni-Flash -Thinking
DailyOmni	69.8(Tang et al., 2025)	72.7	75.8	<b>76.2</b>
VideoHolmes	55.6(Tang et al., 2025)	49.5	<b>57.3</b>	<b>57.3</b>



Thank you