

Image Understanding (re)

- ViT
- CLIP
- ImageBind

Vision Transformer (2021)

- Patchify Image
- Patch Linear Projection
- Image Classification

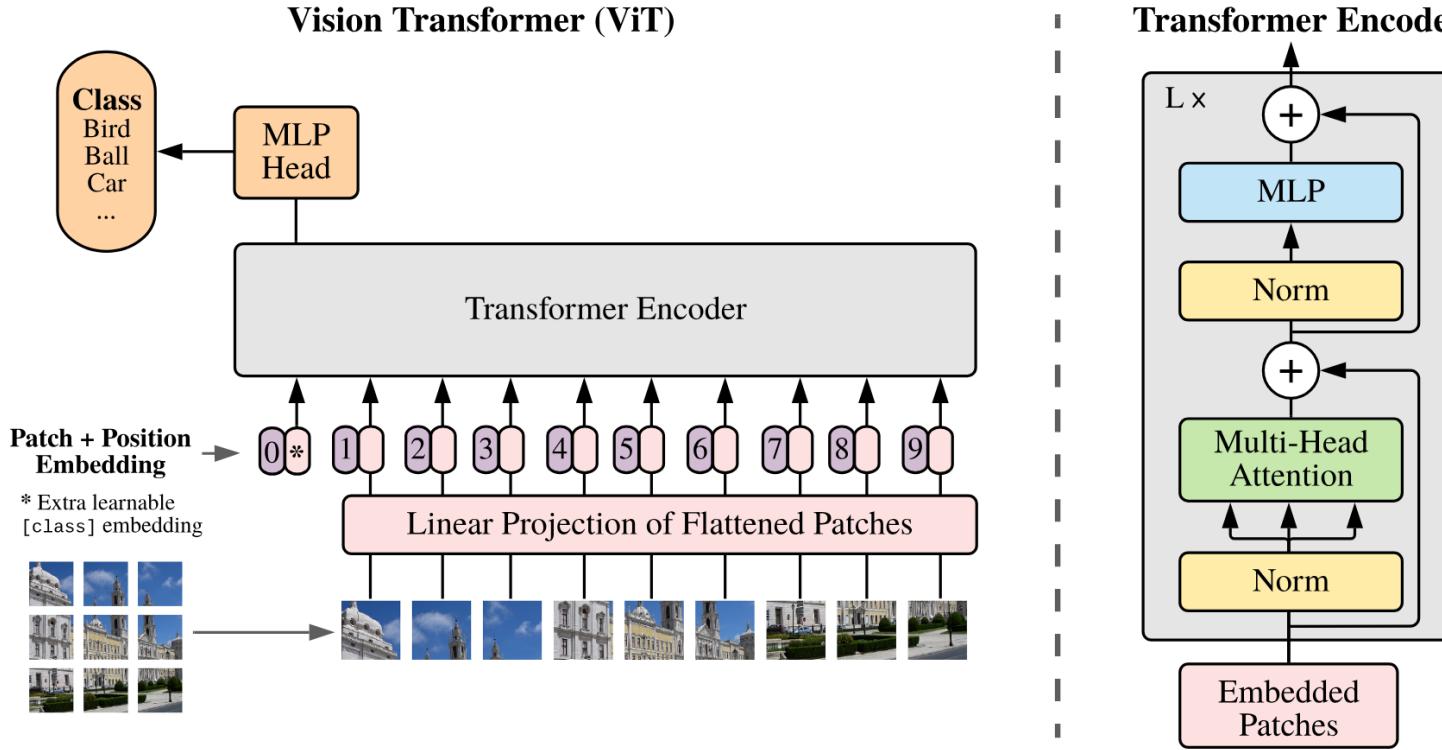
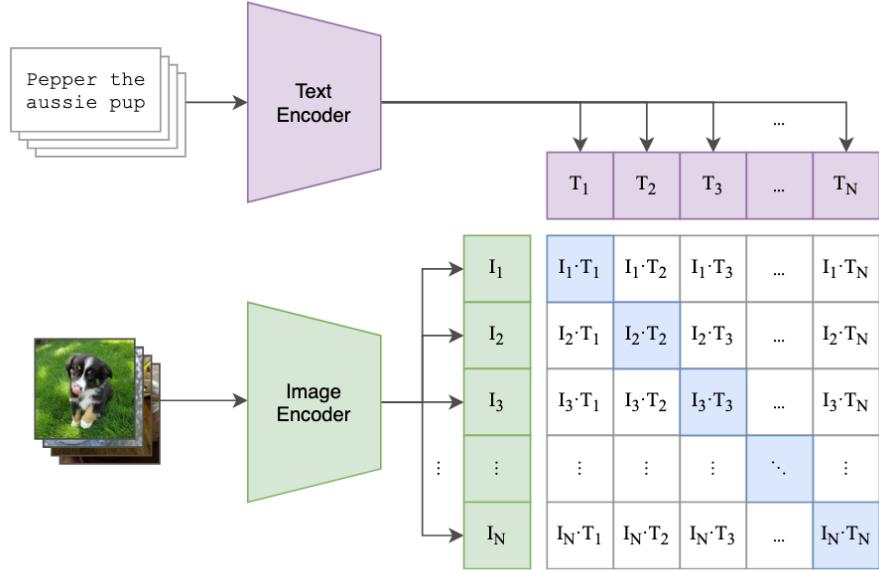


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

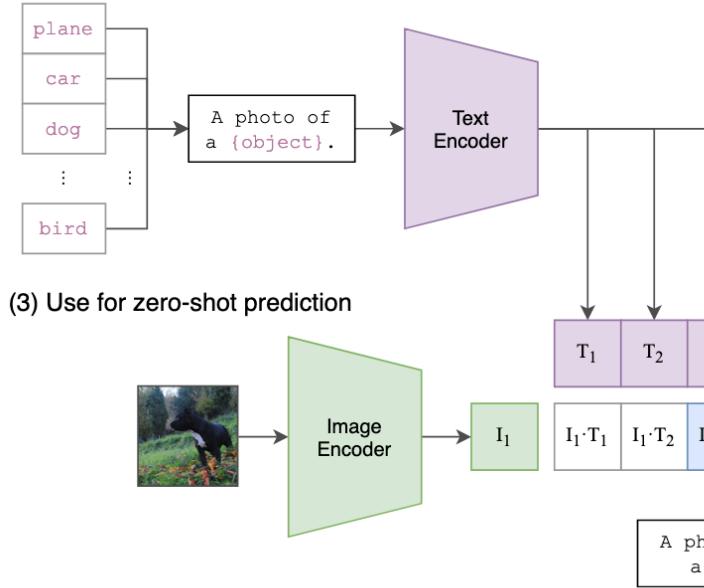
CLIP (2021)

- ViT Image Encoder
- Text Transformer
- One vector for each image and text
- Text-Image Contrastive Loss

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

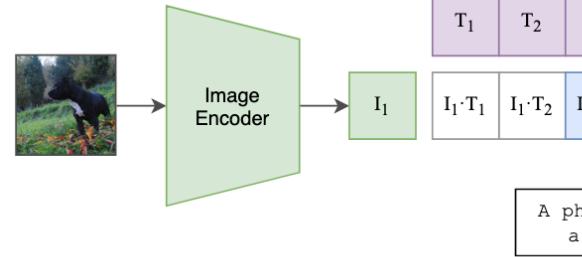


Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier on some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

CLIP (2021)

```
# extract feature representations of each modal
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss = (loss_i + loss_t)/2
```

ImageBind (2023)

- Mostly same as CLIP
- More modalities supported

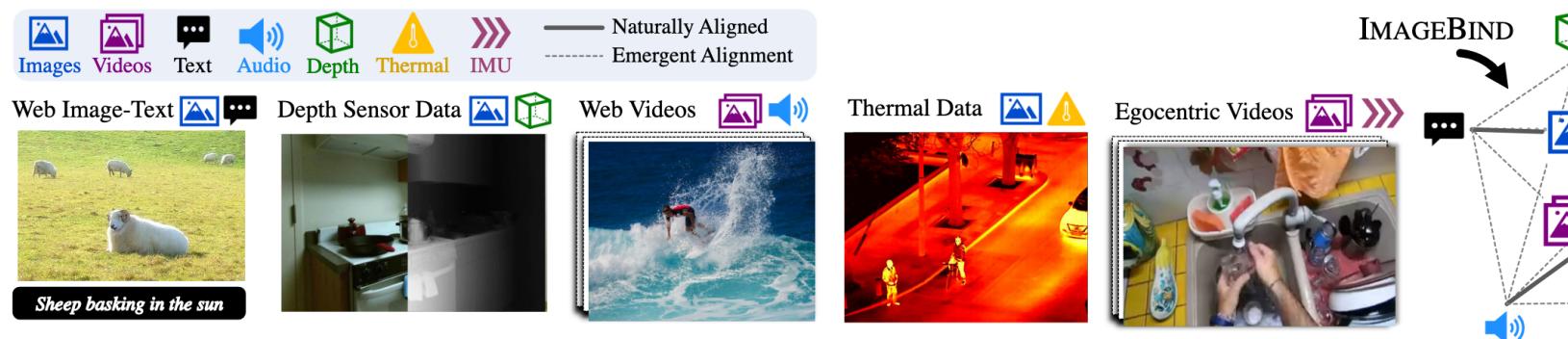


Figure 2. IMAGEBIND overview. Different modalities occur naturally aligned in different data sources, for instance image+video+audio in web data, depth or thermal information with images, IMU data in videos captured with egocentric cameras, IMAGEBIND links all these modalities in a common embedding space, enabling new emergent alignments and capabilities.

Comparison

Feature	Vision Transformer (ViT)	CLIP	Image Text
Training Objective	Classification	Image-text contrastive	Multimodal contrastive
Text Modality	✗	✓	✓
Image Modality	✓	✓	✓
Audio Modality	✗	✗	✓

Image Understanding (re)

-  ViT
-  CLIP
-  ImageBind

Image Understanding Questions

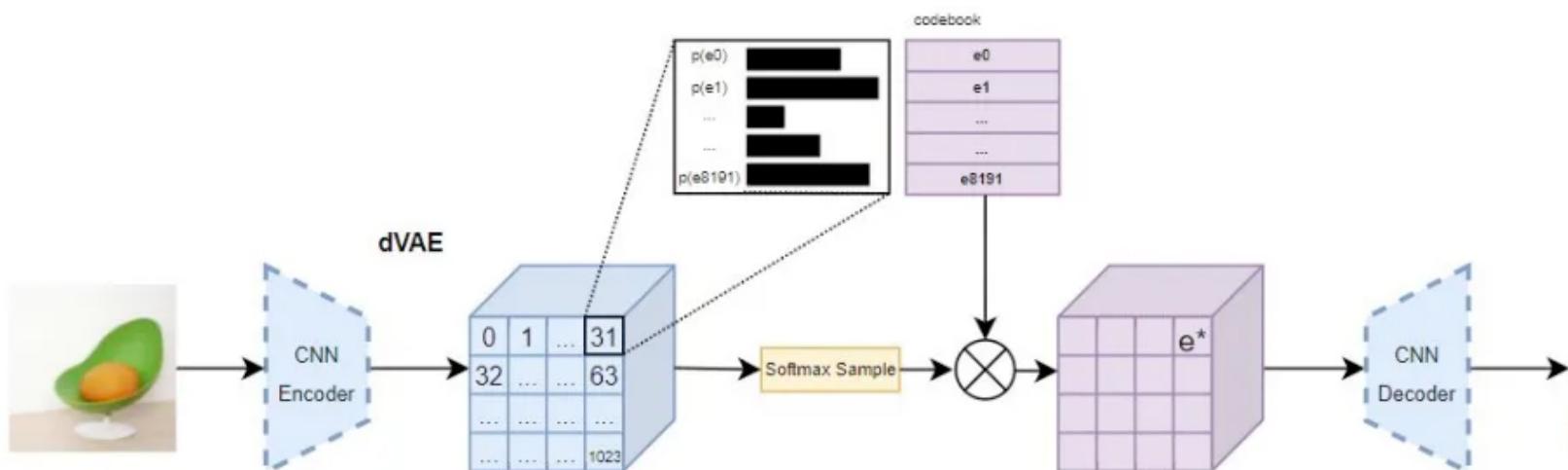
- What advantages does the transformer architecture bring to image processing compared to CNNs?
- How does CLIP's ability to learn from natural language supervision make it more flexible than traditional image classification models ?
- What is shared embedding space for different modalities ?

Image Generation (not LL based)

- Dalle
- Dalle2
- VAR

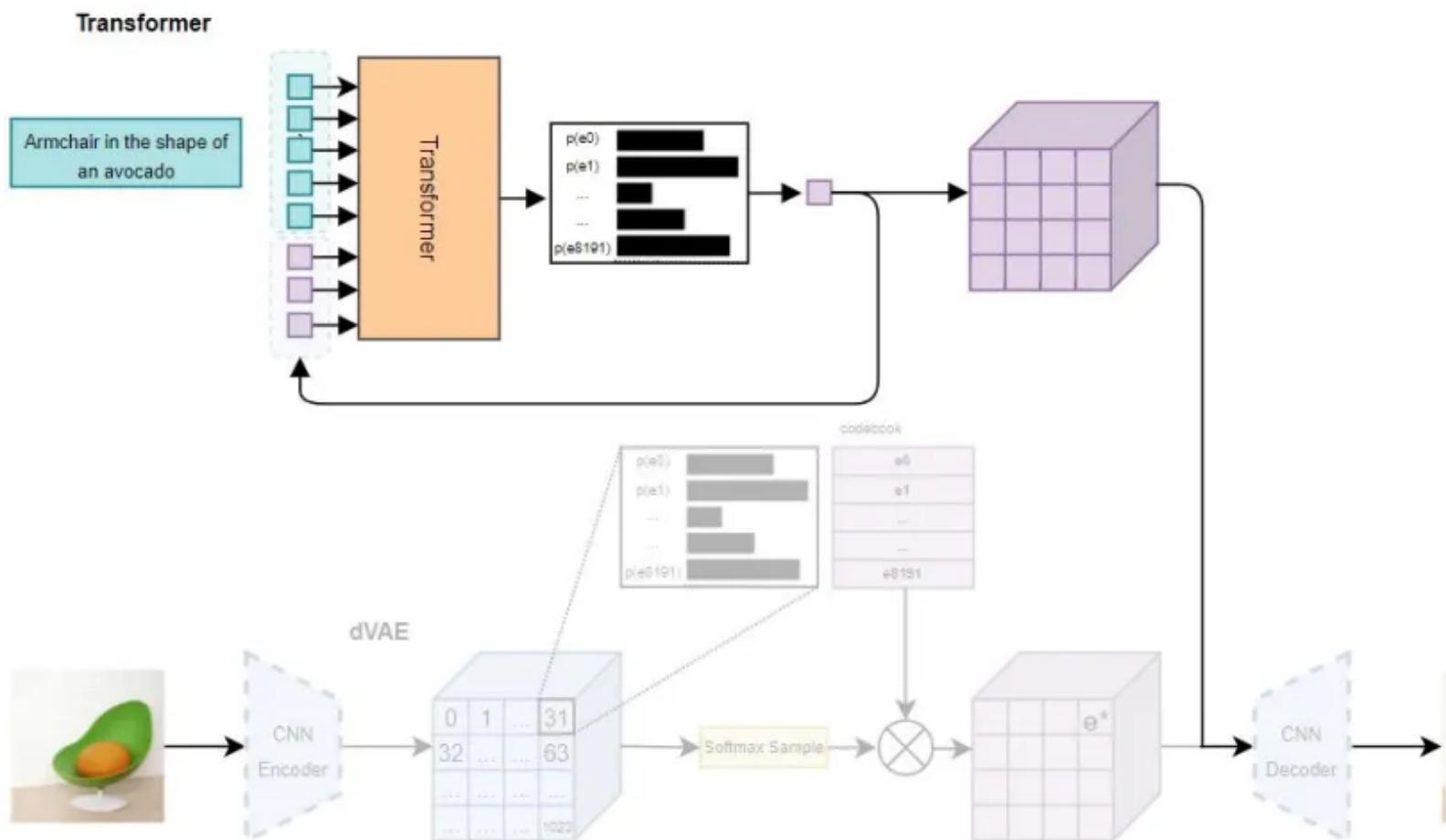
Dalle (2021)

- Compress Image to Discrete Latent Space
- Transformer Causal Modeling Of Text and Im



Dalle (2021)

- Compress Image to Discrete Latent Space
- Transformer Causal Modeling Of Text and Im



Dalle2 (2022)

- Pretrained CLIP (Joint Embedding Space)
- Prior Model for Text to Image embedding Tra
- Diffusion Decoder

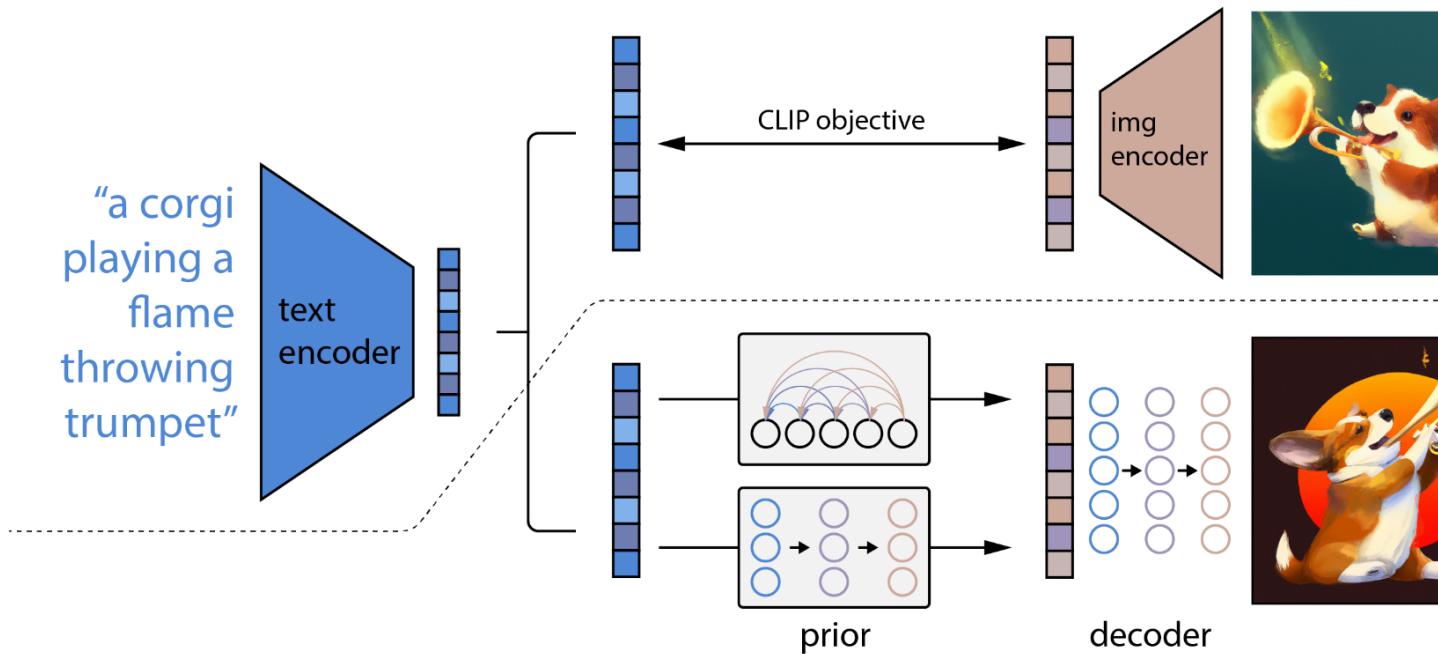


Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process through which we learn a joint representation space for text and images. Below the dotted line, we show the text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or diffusion prior to produce an image embedding, and then this embedding is used to condition a diffusion decoder to produce a final image. Note that the CLIP model is frozen during training of the prior and decoder.

VAR (2024)

- Visual Autoregressive Modeling via Next-Scale Prediction
- VQVAE for Image Compression

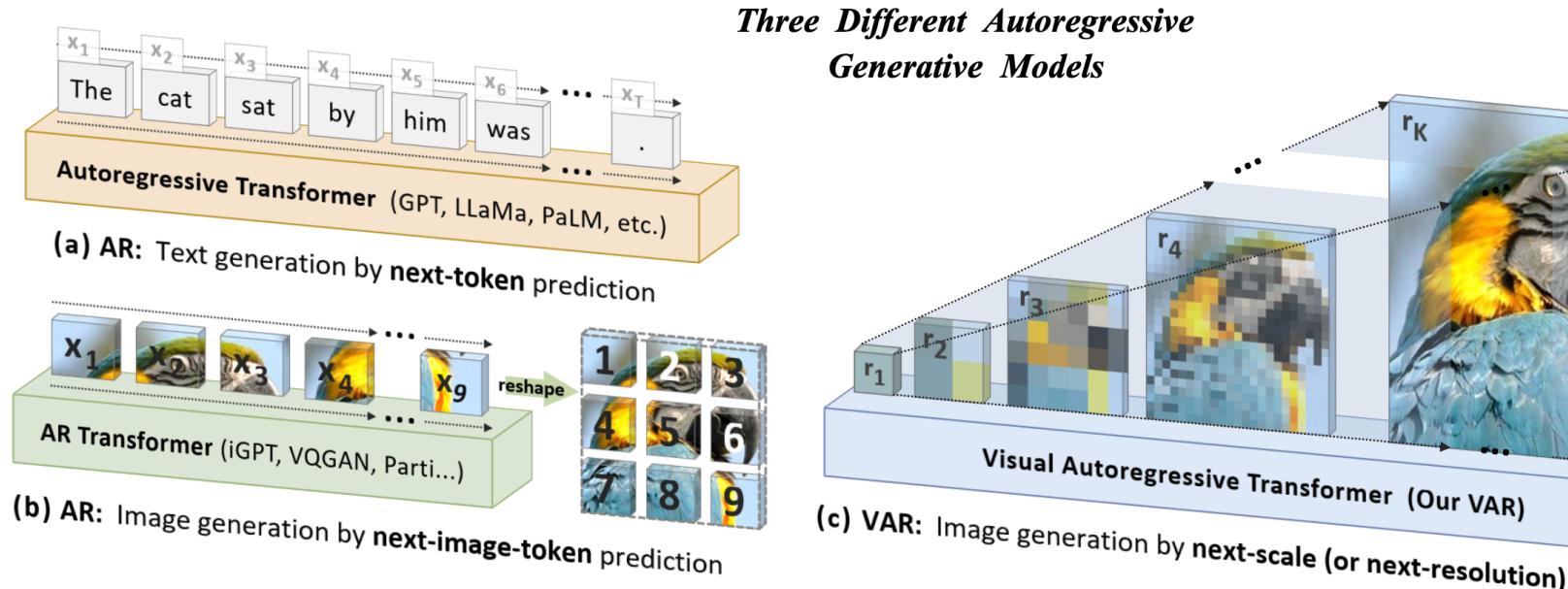


Figure 2: Standard autoregressive modeling (AR) vs. our proposed visual autoregressive model (VAR).

(a) AR applied to language: sequential text token generation from left to right, word by word; (b) AR applied to images: sequential visual token generation in a raster-scan order, from left to right, top to bottom; (c) VAR for images: multi-scale token maps are autoregressively generated from coarse to fine scales (lower resolutions), with parallel token generation within each scale. VAR requires a multi-scale VQVAE to handle the hierarchical nature of images.

Image Generation

-  Dalle - Autoregressive image generation
-  Dalle2 - Diffusion image generation
-  VAR - Visual Autoregressive Modeling via Scale Prediction

Questions

- Dalle{1,2} differences ?
- How to teach LLMs to generate images ?

LLM-based Approaches

Understanding

- Early Fusion - Inputs Concatenation
- Deep Fusion - Cross Attention

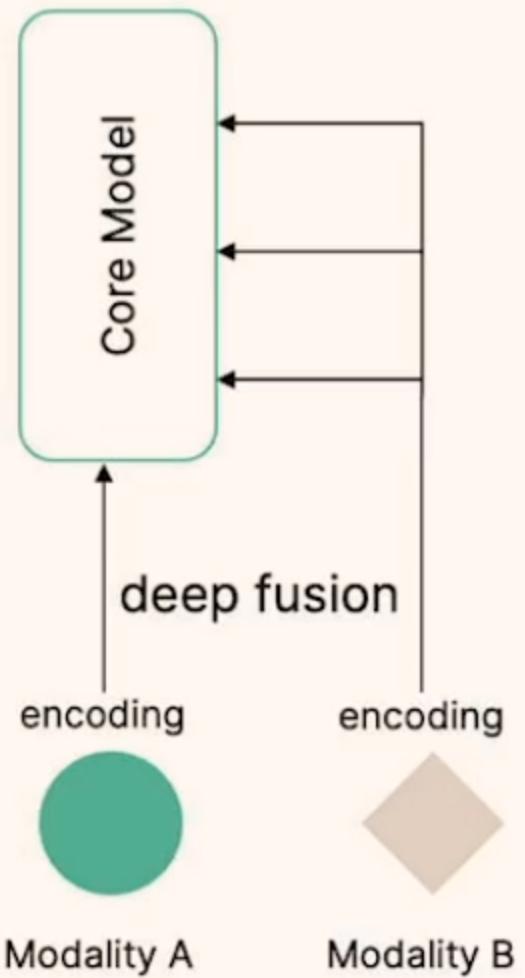
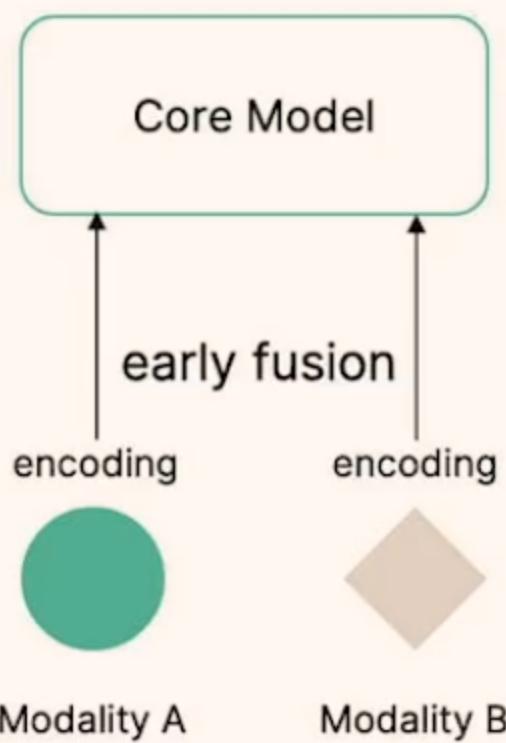


Image Credit

Модели:

- Flamingo
- Formage
- Encodec (2022)
- Llava

Flamingo (2022)

- CLIP Image Encoder
- ! Interleaved Data
- ! Deep Fusion

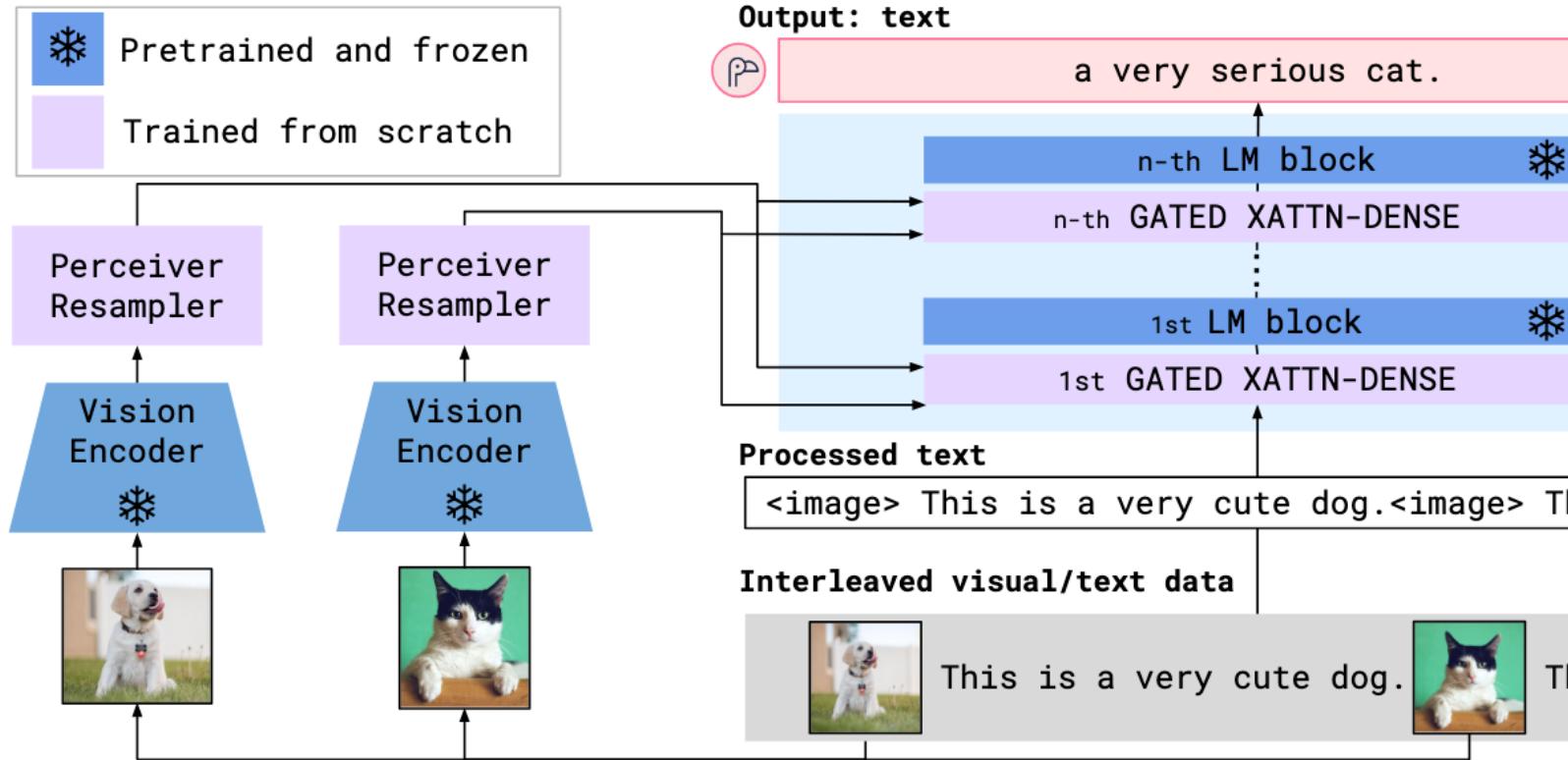
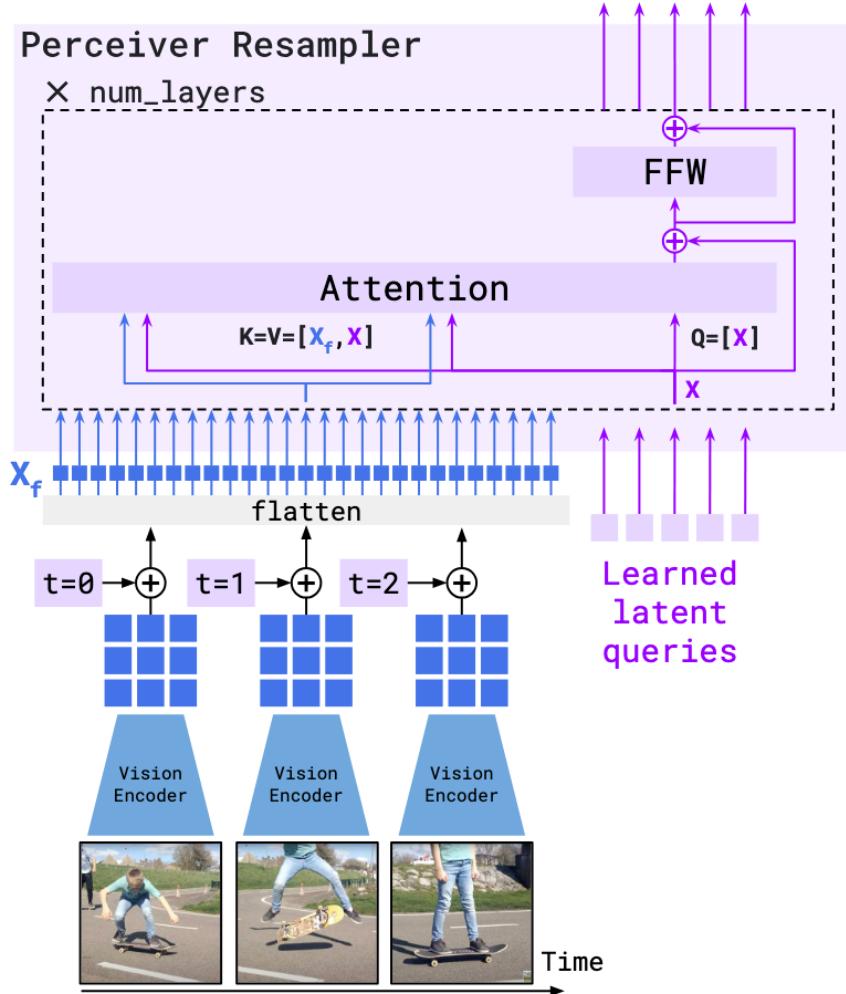


Figure 3: **Flamingo architecture overview.** Flamingo is a family of visual language models that take as input visual data interleaved with text and produce free-form text as output.

Flamingo (2022)

! Perceiver Resampler

- Fixed image embeddings size
- Learned Query Vectors



```

def perceiver_resampler(
    x_f, # The [T, S, d] visual features (T=time, S=spatial)
    time_embeddings, # The [T, 1, d] time pos embeddings
    x, # R learned latents of shape [R, d]
    num_layers, # Number of layers
):
    """The Perceiver Resampler model."""

    # Add the time position embeddings and flatten.
    x_f = x_f + time_embeddings
    x_f = flatten(x_f) # [T, S, d] -> [T * S, d]
    # Apply the Perceiver Resampler layers.
    for i in range(num_layers):
        # Attention.
        x = x + attention_i(q=x, kv=concat([x_f, x]))
        # Feed forward.
        x = x + ffw_i(x)
    return x

```

Figure 5: The Perceiver Resampler module maps a *variable* size grid of spatio-temporal features output by the Vision Encoder to a *fixed* number of output tokens (five in the figure) independently from the input image resolution or the number of input video frames. This transformation has a set of learned latent vectors as queries, and the keys and values are a concatenation of spatio-temporal visual features with the learned latent vectors.

Fromage (2023)

- CLIP Image Encoder
- Opt LLM
- Early Fusion
- ! Image Retrieval

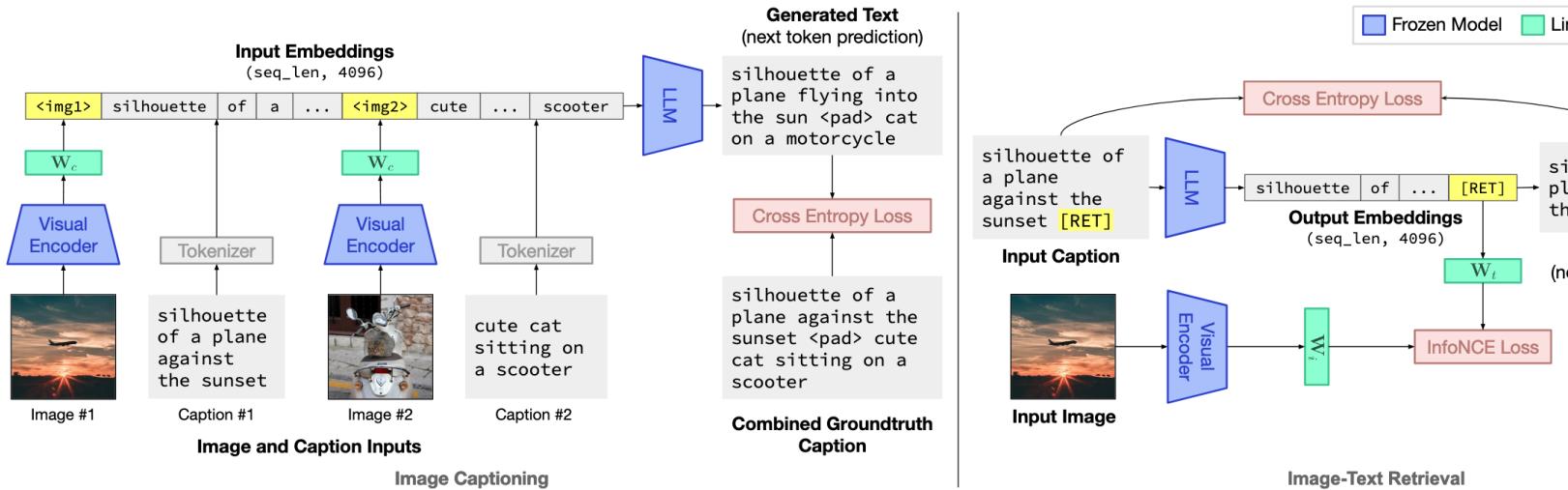


Figure 2. Overview of the FROMAGe architecture. FROMAGe is a model trained on image-text pairs for image captioning and retrieval. It is capable of processing arbitrarily interleaved image and text inputs, and producing interleaved images and te

BLIP2 (2023)

- CLIP Image Encoder
- Opt LLM / Flan T5
- Early Fusion
- ! QFormer for Image Embedding
- Freezed LLM

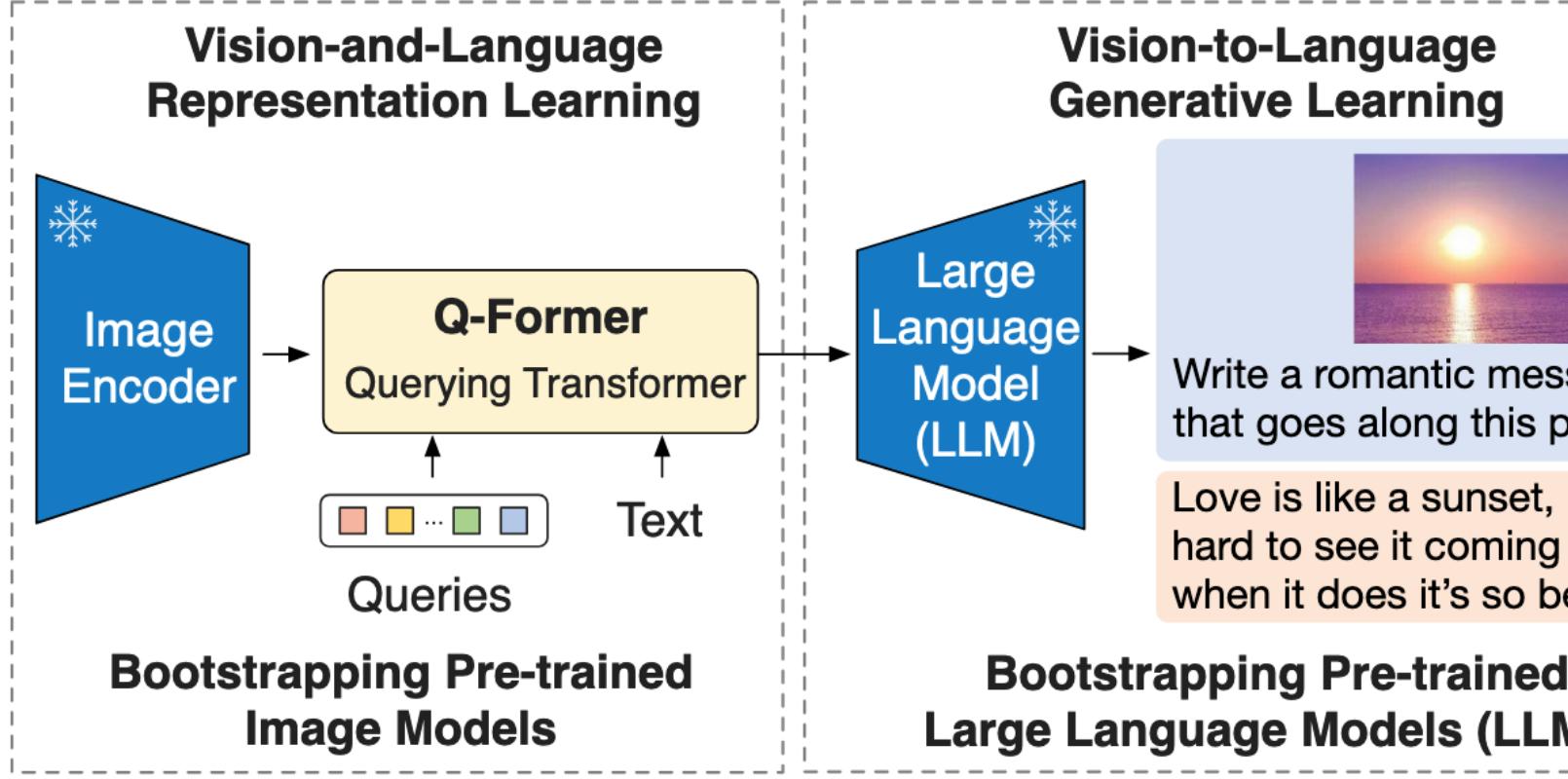


Figure 1. Overview of BLIP-2’s framework. We pre-train a lightweight Querying Transformer following a two-stage strategy to bridge the modality gap. The first stage bootstraps vision-and-language representation learning from a frozen image encoder. The second stage bootstraps vision-to-language generative learning from a frozen LLM, which enables zero-shot instructed image-text generation (see Figure 4 for more examples).

Llava (2023)

- CLIP Image Encoder
- Vicuna-7B LLM
- 2-stage training:
 - Feature Alignment
 - End-to-End SFT

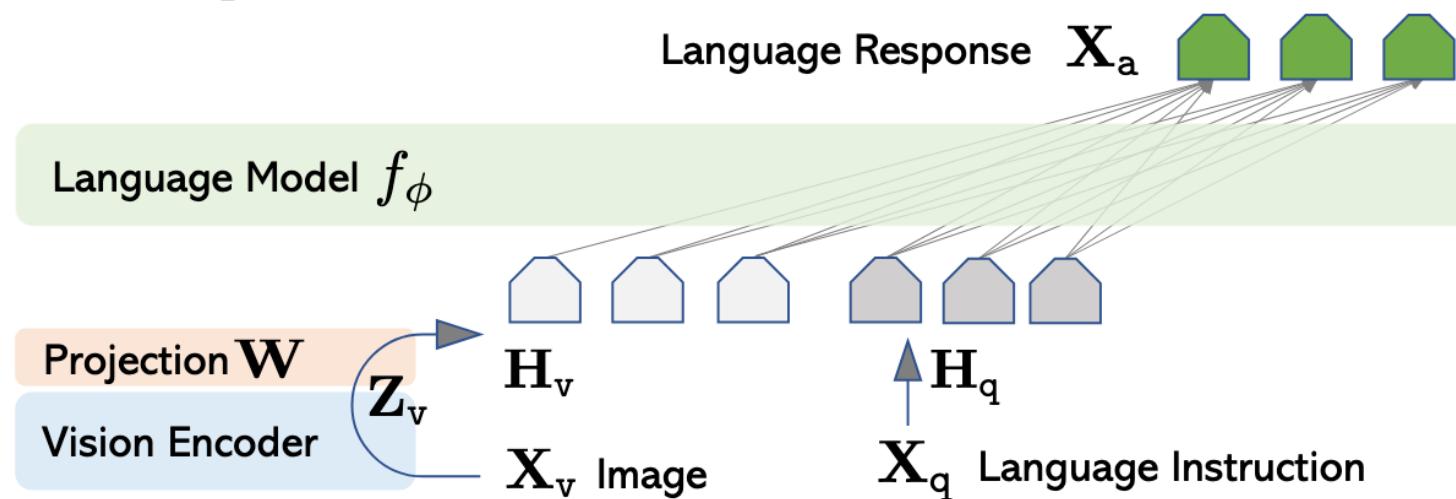


Figure 1: LLaVA network architecture.

LLM Image Understanding Comparison

Feature / Model	Flamingo	Fromage	BLIP2	LLM
Fusion Type	Deep	Early	Early	Early
Interleaved Data	✓	✗	✗	✗
Image Retrieval	✗	✓	✓	✗
Instruction Tuning	✗	✗	✓	✓

LLM Image Understanding Questions

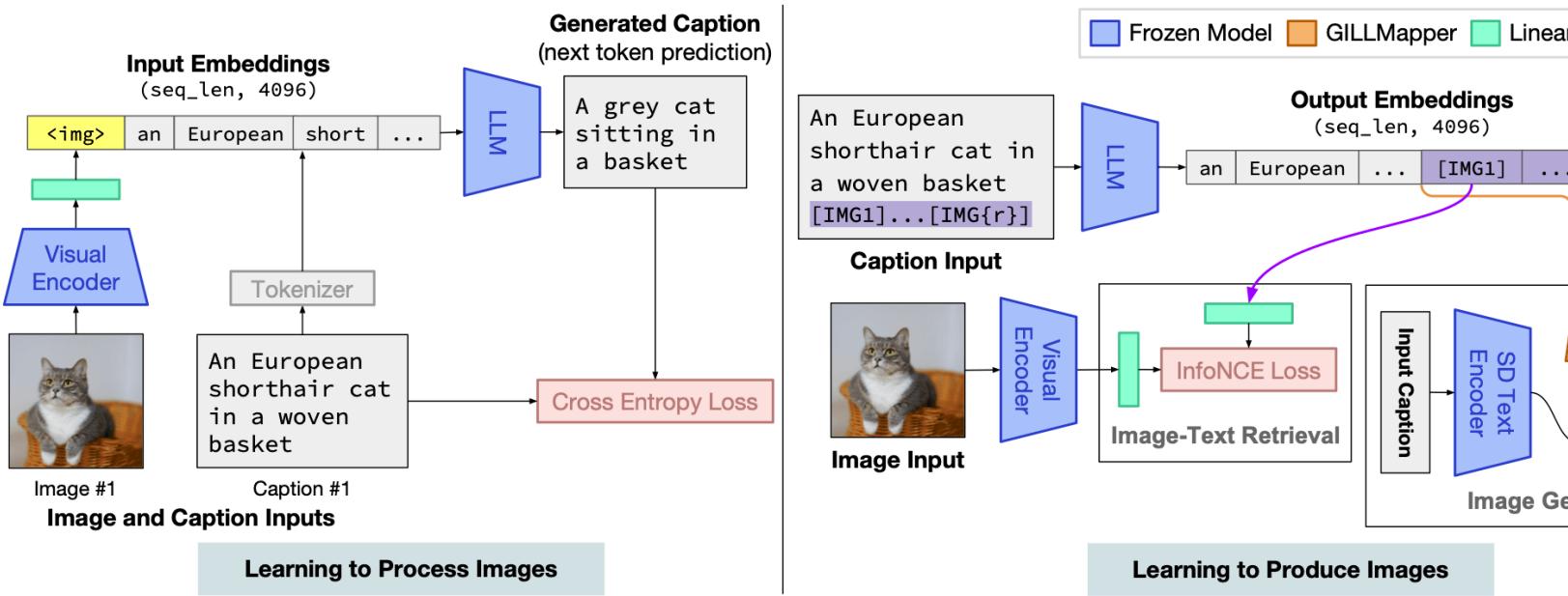
- Why **CLIP ViT** is utilized in most MLLM models
- DeepFusion and Encoder-Decoder Transformer
Difference ?
- Flamingo Perceiver Resampler vs Blip2 QFormer

LLM Image Generation / Editing

- GILL
- Mgie
- Bagel

GILL (2023)

- Image Understanding
- Image Retrieval and Generation
- LLM learns to generate image embeddings for generation model
- 8 visual tokens



GILL (2023)

- How GillMapper looks like ?

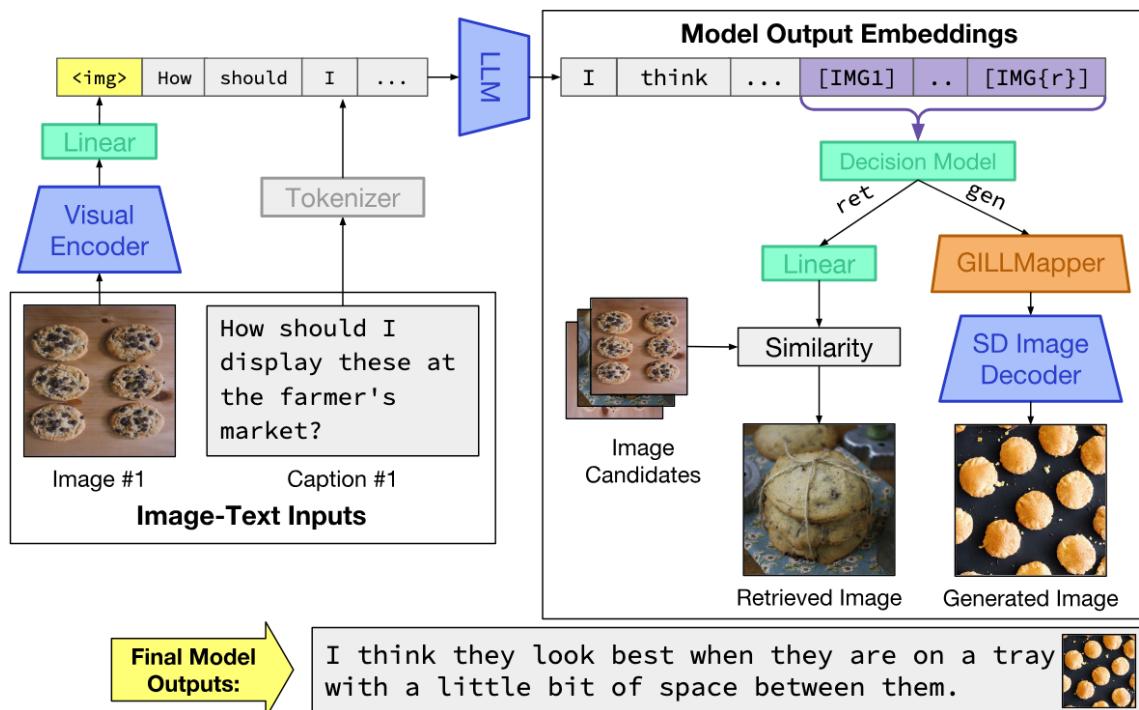


Figure 3: Inference time procedure for GILL. The model takes in image and text inputs, and produces text interleaved with image embeddings. After deciding whether to retrieve or generate for a particular set of tokens, it returns the appropriate image outputs.

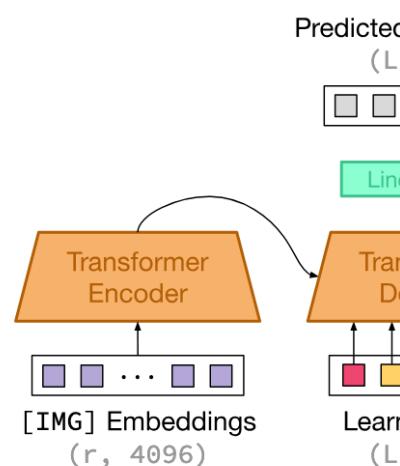
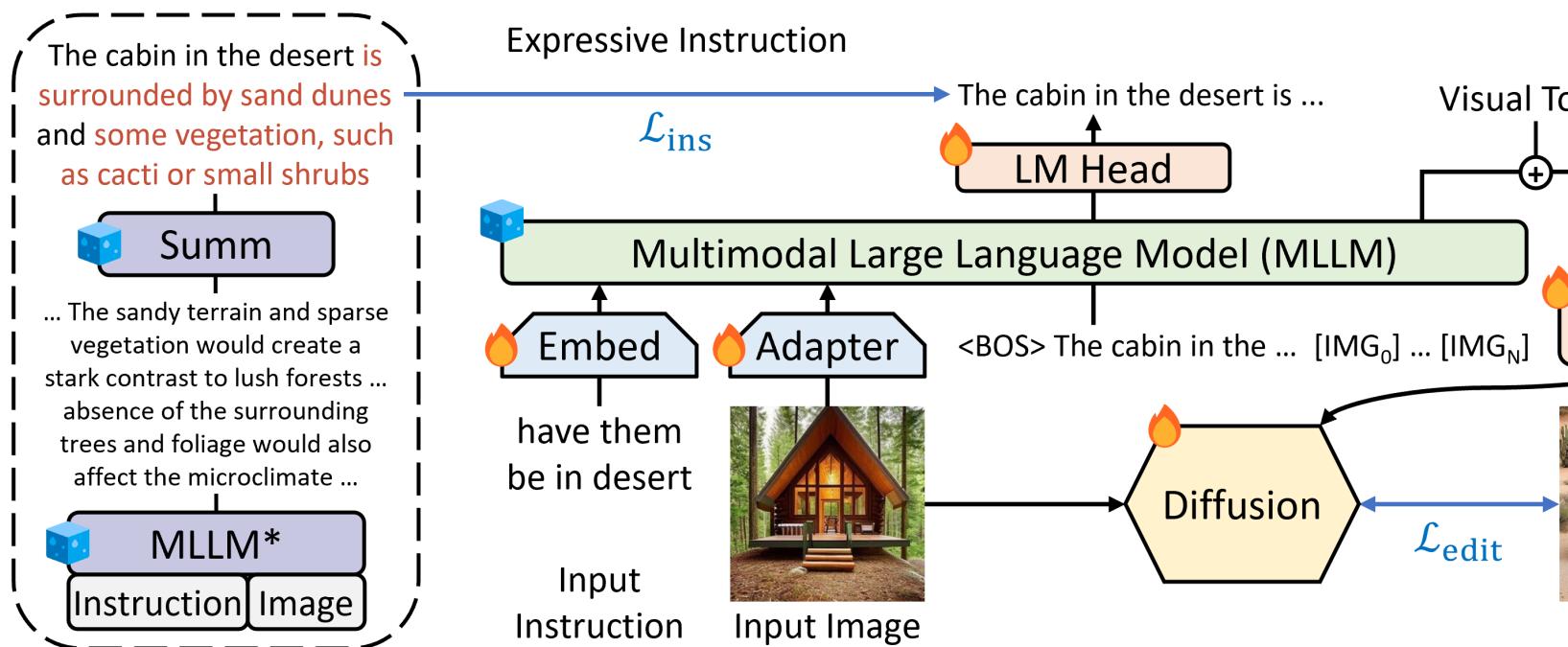


Figure 4: GILLMapper architecture. It is conditioned on the hidden [IMG] embeddings and a sequence of query embedding vectors.

Mgie (2024) Apple

- ! Task: Instruction-based editing
- Expressive Instructions
- Mostly inspired by GILL



Bagel (2025) ByteDance

- ! Tasks: Understanding, Editing, Style Transfer
- Thinking before Generation
- Data
- Qwen2.5 LLM
- SigLIP2 Image Encoder
- FLUX VAE

Language Response



Next Token Prediction

FFN

Image/Multi-Image/Video



Velocity Prediction

FFN

Multi-modal Self Attention

QKV

QKV

Und. Expert

Gen. Expert



Text Tokenizer

Und Encoder

Gen Encoder

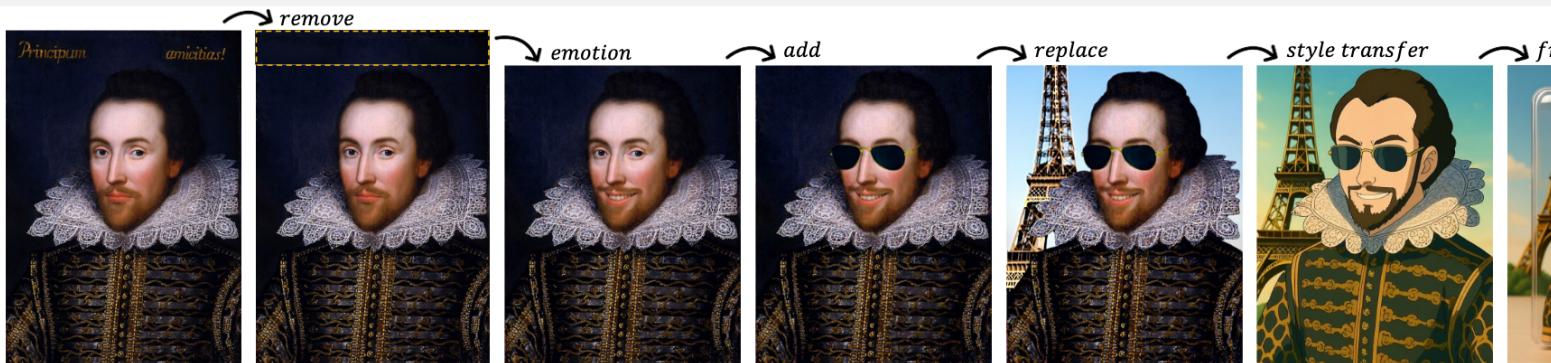
Bagel (2025) ByteDance

Samples

Image Generation



Image Editing



Free-form Manipulation



LLM Image Generation / Editing

-  GILL
-  Mgie
-  Bagel

LLM Image Generation / Editing Metrics

- Side-by-side User Study

Materials

- Елизавета Гончарова | Мультимодальные подходы и LLM
- Максим Куркин | MLLMs
- Обзор на русском от Максима мультимодальных подходов
- Molmo - VLLM от allenai с открытым кодом тренировки
- InternVL-2.5 - SOTA VLLM Image Encoder
- Malvina - как решают задачу pixel-perfect редактирования изображений в команде Gigachat

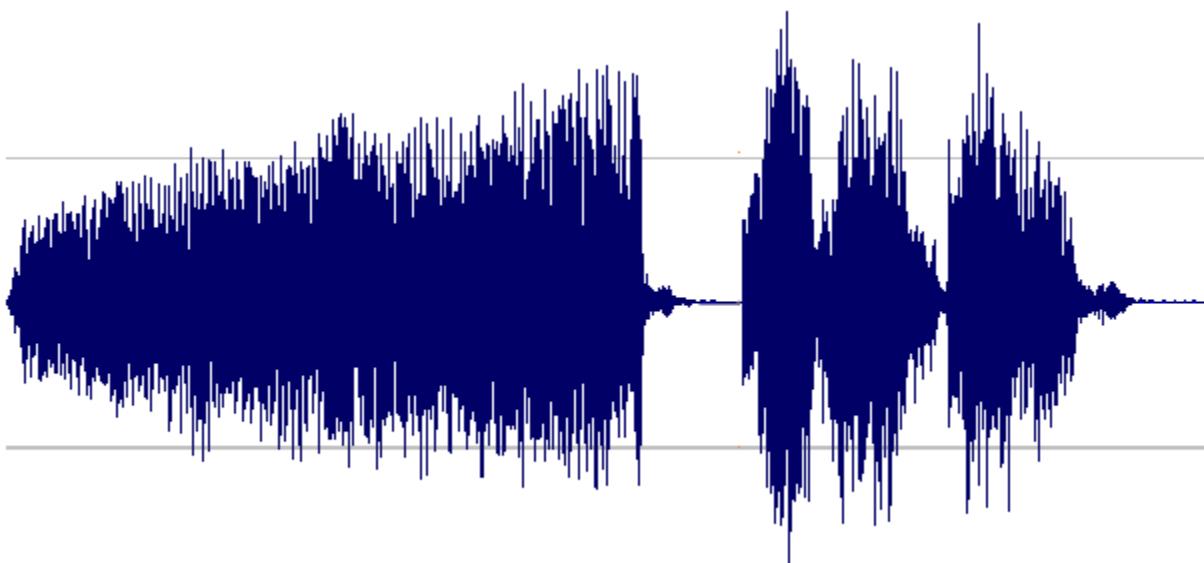
Audio Modality

Audio Representations

- Waveform
- Mel-Spectrogram

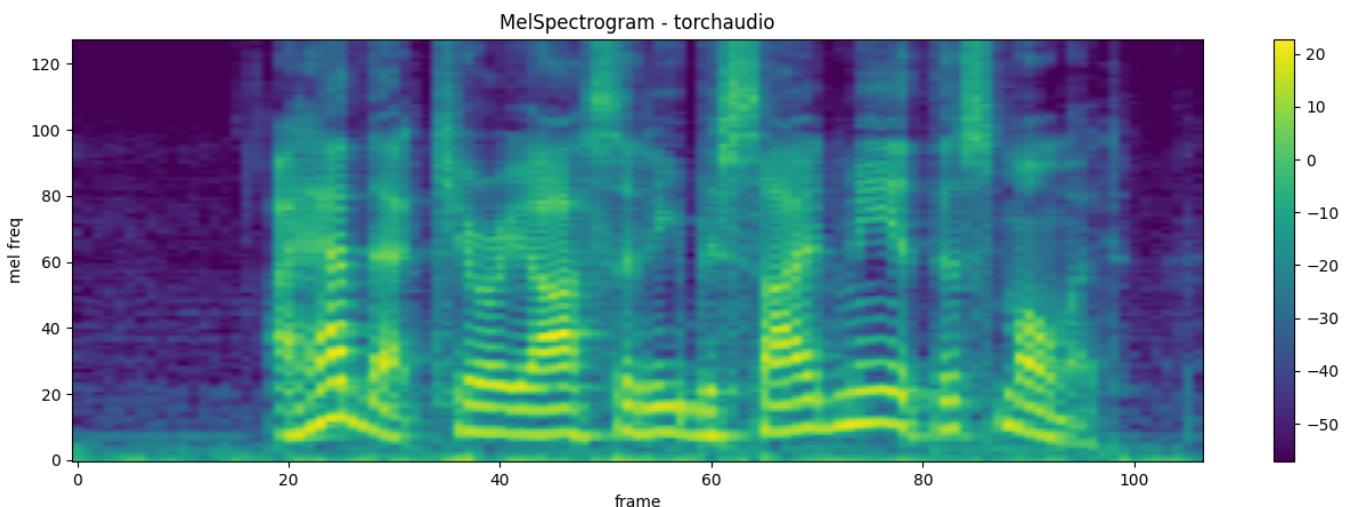
Waveform

- Raw audio samples
- Time series
- Sample rate: 16kHz-48kHz



Mel-Spectrogram

- Semantically compressed (but approximately same size in bytes as waveform)
- Windowed Fourier Transform
- Hop size = 10 ms (Time frames)
- Mel bins = 40 to 128 (Frequency bins)



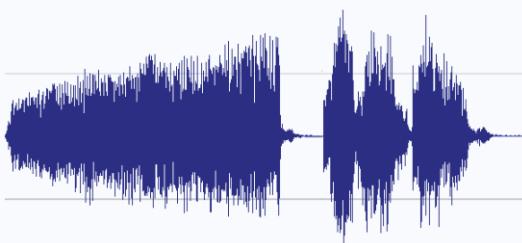
Audio Codecs



Goal:

Discrete Audio Representation for Causal Mod
(and/or audio compression)

Waveform (16kHz)



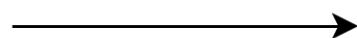
Discrete Codes (<300Hz)

44	200	10	5	19	76
----	-----	----	---	----	----	------

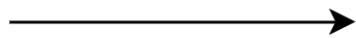
Waveform (16kHz)



Compression



Reconstruction



Audio Codecs

 Models:

- Encodec (2022)
- Mimi (2024)

Encodec (2022)

Architecture

- Varying Bitrates (Residual Vector Quantization)

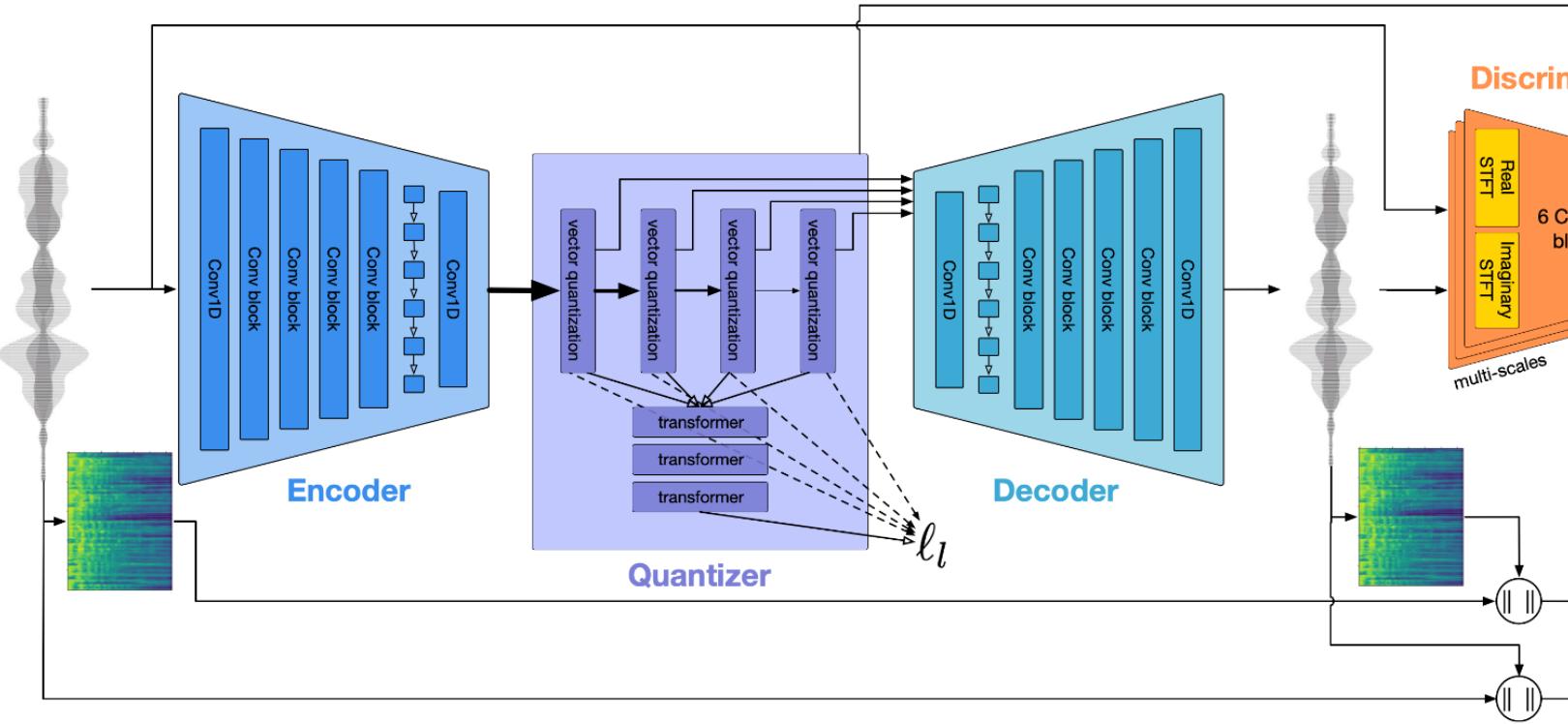


Figure 1: ENCODEC : an encoder decoder codec architecture which is trained with reconstruction loss (ℓ_t) as well as adversarial losses (ℓ_g for the generator and ℓ_d for the discriminator). The residual quantization commitment loss (ℓ_w) applies only to the encoder. Optionally, we train a small Transformer language model for entropy coding over the quantized units with ℓ_l , which reduces bandwidth even further.

Encodec (2022)

LOSS

Reconstruction loss

Reconstruction Loss. The reconstruction loss term is comprised of a time and a frequency domain term. We minimize the L1 distance between the target and compressed audio over the time domain term $\ell_t(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|_1$. For the frequency domain, we use a linear combination between the L1 and L2 losses over the mel-spectrogram using several time scales (Yamamoto et al., 2020b; Gritsenko et al., 2020).

$$\ell_f(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{|\alpha| \cdot |s|} \sum_{\alpha_i \in \alpha} \sum_{i \in e} \|\mathcal{S}_i(\mathbf{x}) - \mathcal{S}_i(\hat{\mathbf{x}})\|_1 + \alpha_i \|\mathcal{S}_i(\mathbf{x}) - \mathcal{S}_i(\hat{\mathbf{x}})\|_2,$$

Encoder (2022)

Loss

Discriminative Loss

The adversarial loss for the generator is constructed as follows, $\ell_g(\hat{\mathbf{x}}) = \frac{1}{K} \sum_k \max(0, 1 - D_k(\hat{\mathbf{x}}))$, where K is the number of discriminators. Similarly to previous work on neural vocoders (Kumar et al., 2019; Karras et al., 2020; You et al., 2021), we additionally include a relative feature matching loss for the generator.

$$\ell_{feat}(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{KL} \sum_{k=1}^K \sum_{l=1}^L \frac{\|D_k^l(\mathbf{x}) - D_k^l(\hat{\mathbf{x}})\|_1}{\text{mean}(\|D_k^l(\mathbf{x})\|_1)},$$

where the mean is computed over all dimensions, (D_k) are the discriminators, and L is the number of layers per discriminator. The discriminators are trained to minimize the following hinge-loss adversarial loss: $L_d(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{K} \sum_{k=1}^K \max(0, 1 - D_k(\mathbf{x})) + \max(0, 1 + D_k(\hat{\mathbf{x}}))$, where K is the number of discriminators. To prevent that the discriminator tend to overpower easily the decoder, we update its weight with a probability of 0.8 at 24 kHz, and 0.5 at 48 kHz.

Encoder (2022)

Loss

RVQ Commitment Loss

VQ commitment loss. As mentioned in Section 3.2, we add a commitment loss l_w between the encoder, and its quantized value, with no gradient being computed for the quantized value residual step $c \in \{1, \dots, C\}$ (with C depending on the bandwidth target for the current batch), now current residual and $q_c(\mathbf{z}_c)$ the nearest entry in the corresponding codebook, we define l_w as

$$l_w = \sum_{c=1}^C \|\mathbf{z}_c - q_c(\mathbf{z}_c)\|_2^2.$$

Overall, the generator is trained to optimize the following loss, summed over the batch,

$$L_G = \lambda_t \cdot \ell_t(\mathbf{x}, \hat{\mathbf{x}}) + \lambda_f \cdot \ell_f(\mathbf{x}, \hat{\mathbf{x}}) + \lambda_g \cdot \ell_g(\hat{\mathbf{x}}) + \lambda_{feat} \cdot \ell_{feat}(\mathbf{x}, \hat{\mathbf{x}}) + \lambda_w \cdot \ell_w(w),$$

where λ_t , λ_f , λ_g , λ_{feat} , and λ_w the scalar coefficients to balance between the terms.

Encodec (2022)

Loss

Balancer

Balancer. We introduce a loss balancer in order to stabilize training, in particular the varying gradients coming from the discriminators. We also find that the balancer makes it easier to reason about different loss weights, independently of their scale. Let us take a number of losses $(\ell_i)_i$ that depend on the output of the model \hat{x} . We define $g_i = \frac{\partial \ell_i}{\partial \hat{x}}$, and $\langle \|g_i\|_2 \rangle_\beta$ the exponential moving average of the last training batches. Given a set of weights (λ_i) and a reference norm R , we define

$$\tilde{g}_i = R \frac{\lambda_i}{\sum_j \lambda_j} \cdot \frac{g_i}{\langle \|g_i\|_2 \rangle_\beta}.$$

We then backpropagate into the network $\sum_i \tilde{g}_i$, instead of the original $\sum_i \lambda_i g_i$. This changes the optimization problem but allows to make the λ_i interpretable irrespectively of the natural scale of each loss. If $\lambda_i = 1$ then each weight can be interpreted as the fraction of the model gradient that come from the corresponding loss. We take $R = 1$ and $\beta = 0.999$. All the generator losses from Eq. (4) fit into the balancer, except the commitment loss, as it is not defined with respect to the output of the model.

Encodec (2022)

Evaluation

Table 1: MUSHRA scores for Opus, EVS, Lyra-v2, and ENCODEC for various bandwidths in streamable setting. Results are reported across different audio categories (clean speech, noisy speech, music), sampled at 24 kHz. We report mean scores and 95% confidence intervals. For ENCODEC, we report the average bandwidth after using the entropy coding described in Section 3.3.

Model	Bandwidth	Entropy Coded	Clean Speech	Noisy Speech	Music Set-1	Music Set-2
Reference	-	-	95.5±1.6	93.9±1.8	93.2±2.5	97.0±1.5
Opus	6.0 kbps	-	30.1±2.8	19.1±5.9	20.6±5.8	17.5±2.5
Opus	12.0 kbps	-	76.5±2.3	61.9±2.1	77.8±3.2	65.0±2.5
EVS	9.6 kbps	-	84.4±2.5	80.0±2.4	89.9±2.3	87.5±2.5
Lyra-v2	3.0 kbps	-	53.1±1.9	52.0±4.7	69.3±3.3	42.0±2.5
Lyra-v2	6.0 kbps	-	66.2±2.9	59.9±3.3	75.7±2.6	48.0±2.5
ENCODEC	1.5 kbps	0.9 kbps	49.2±2.4	41.3±3.6	68.2±2.2	66.0±2.5
ENCODEC	3.0 kbps	1.9 kbps	67.0±1.5	62.5±2.3	89.6±3.1	87.0±2.5
ENCODEC	6.0 kbps	4.1 kbps	83.1±2.7	69.4±2.3	92.9±1.8	91.0±2.5
ENCODEC	12.0 kbps	8.9 kbps	90.6±2.6	80.1±2.5	91.8±2.5	92.0±2.5

Audio Codecs

 Models:

-  Encodec (2022)
- Mimi (2024)

Mimi (2024)

Architecture

- NEW Semantic Features Distillation

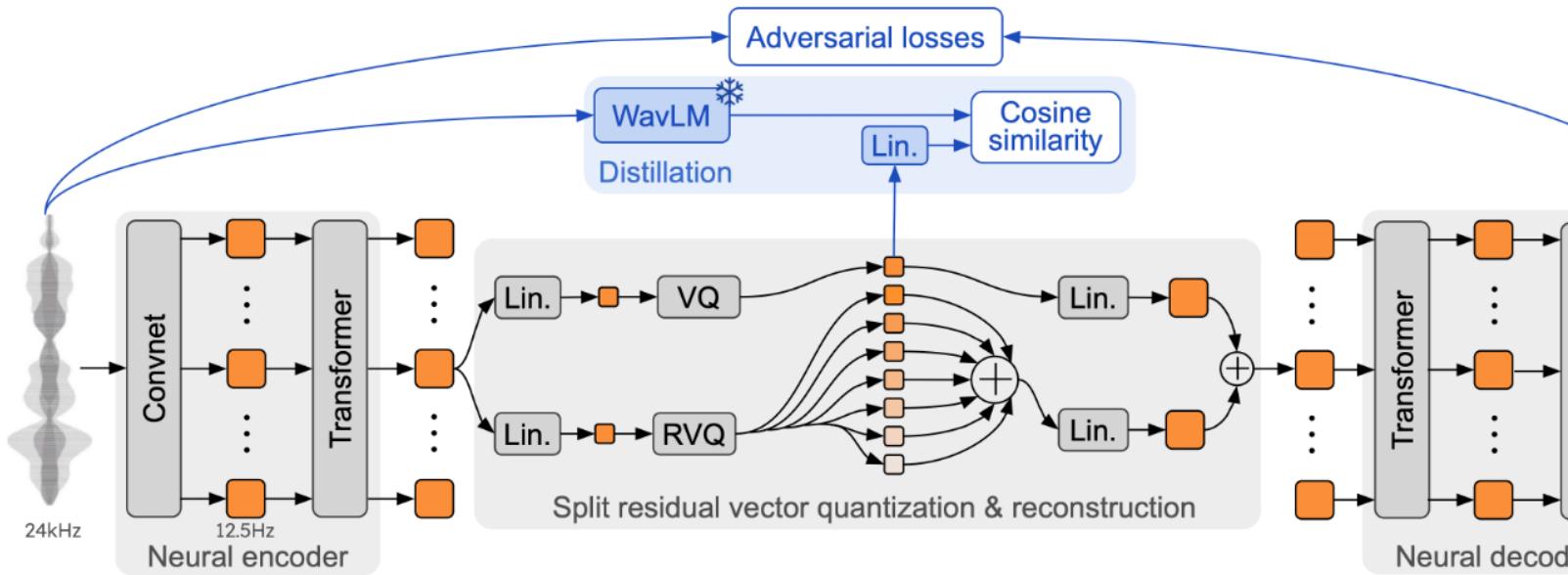


Figure 2: **Architecture and training of Mimi, our neural audio codec, split residual vector quantization.** During training (blue part, top), we dis-

Mimi (2024)

Metrics

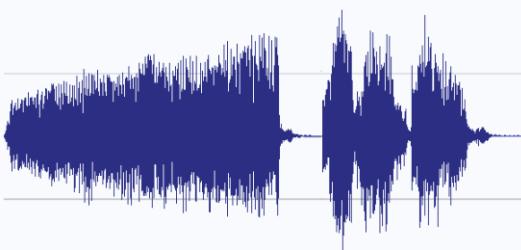
Table 4: **Audio quality evaluation.** Objective and subjective (MUSHRA) eval audio quality for baseline neural audio codecs—RVQGAN (Kumar et al., 2024) tiCodec (Liu et al., 2024), and SpeechTokenizer (Zhang et al., 2024b)— and important variants of Mimi. For a fair comparison with SemantiCodec and Sp enizer, we also include a downsampled version of our codec in the MUSHRA stu the audio sample rate and f_r the codec frame rate. Both Mimi codecs are tra distillation, and either with the same combination of reconstruction and adversar as Encodec (see Section 3.3) or adversarial losses only.

Model	f_s	f_r	bitrate	causal	ABX (\downarrow)	VisQOL (\uparrow)	MOSNet (\uparrow)	MU
Ground Truth	24kHz	-	-	-	-	-	3.08	
RVQGAN	24kHz	75Hz	1.5kbps		-	1.74	2.74	
SemantiCodec	16kHz	50Hz	1.3kbps		42.2%	2.43	3.12	
SpeechTokenizer	16kHz	50Hz	1.5kbps		3.3%	1.53	2.67	
SpeechTokenizer	16kHz	50Hz	4.0kbps		3.3%	3.07	3.10	
Mimi, adv. loss only	24kHz	12.5Hz	1.1kbps	✓	8.7%	1.84	3.10	
Same, downsampled at 16kHz	16kHz	12.5Hz	1.1kbps	✓	-	-	-	
Mimi, non adv. only	24kHz	12.5Hz	1.1kbps	✓	8.1%	2.82	2.89	

Audio Codecs

-  Encoder
-  Mimi

Waveform (16kHz)



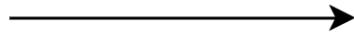
Discrete Codes (<300Hz)

44 | 200 | 10 | 5 | 19 | 76 |

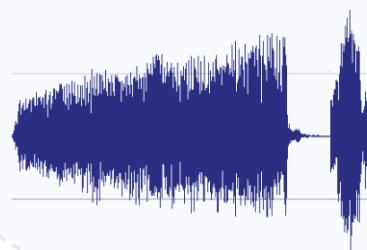
Compression



Reconstruction



Waveform (16kHz)



Audio Codecs

Feature	Mimi	Encodec (Meta)
Input audio sample rate	16 kHz	24,48 kHz
Codec sample rate	12.5 Hz	150-2400 Hz
Varying bandwidth	✗	✓
Semantic distillation	✓	✗

Conversational LLMs

 Audio Understanding

 Audio Generation

 Emotional

 Interruptions handling

 Instructions following (world model)

Conversational LLMs

Live demo of GPT4-o voice variation



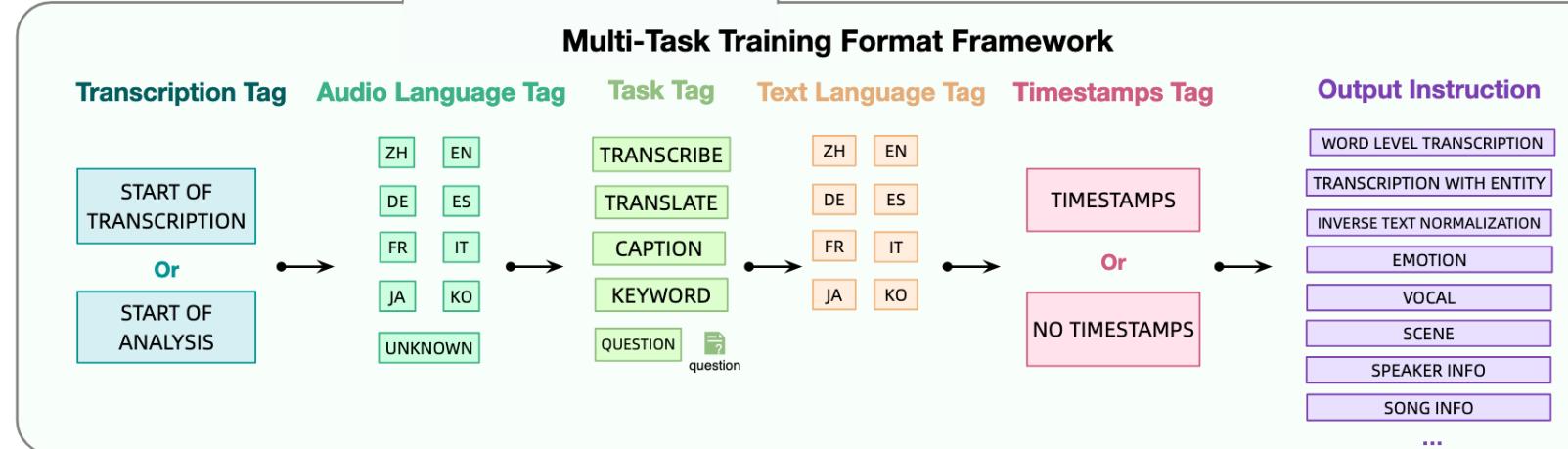
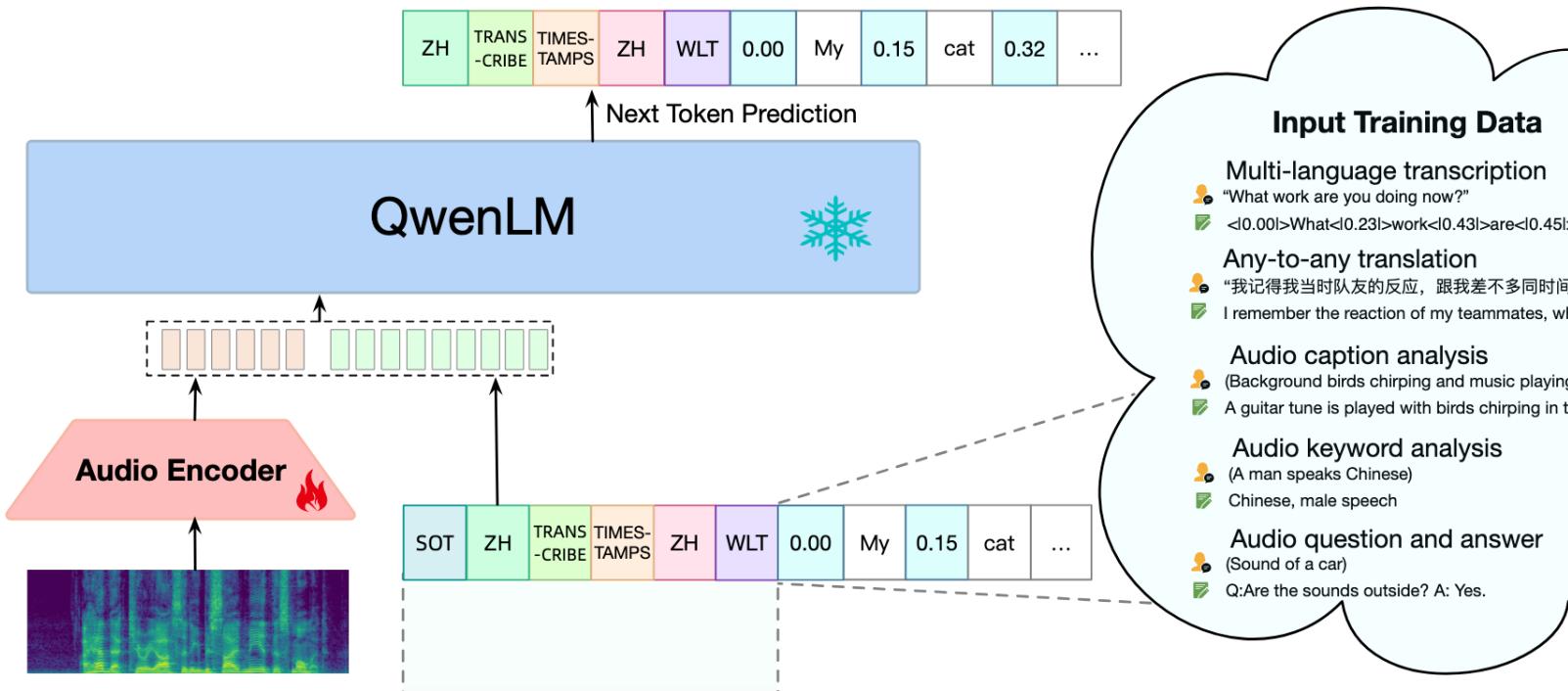
Audio LLMs

- Qwen2.5-Audio ➡️
- Qwen2.5-Omni
- Moshi

Qwen2.5-Audio.

Architecture.

- Whisper-Initialized Audio Encoder
- Qwen2.5-7B LLM Backbone



Qwen2.5-Audio.

Training.

2 stages training:

1. Multitask pretraining:

 Audio Encoder.  LLM.

2. SFT:

 Audio Encoder.  LLM.

Qwen2.5-Audio.

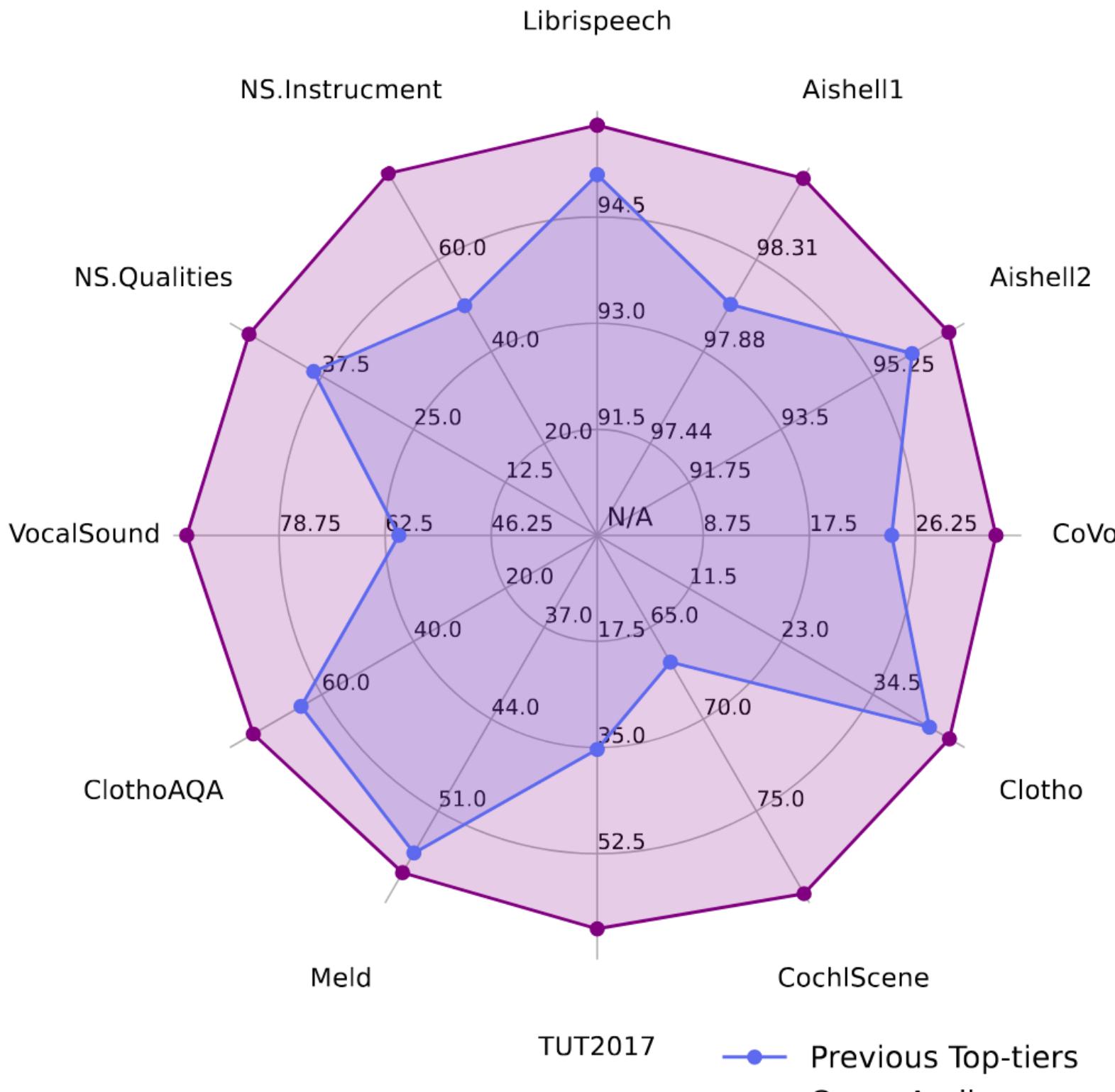
Tasks.

Types	Task	Description
Speech	ASR	Automatic speech recognition (multiple languages)
	S2TT	Speech-to-text translation
	OSR	Overlapped speech recognition
	Dialect ASR	Automatic dialect speech recognition
	SRWT	English speech recognition with word-level timestamps
		Mandarin speech recognition with word-level timestamps
	DID	Dialect identification
	LID	Spoken language identification
	SGC	Speaker gender recognition (biologically)
	ER	Emotion recognition
	SV	Speaker verification
	SD	Speaker diarization
	SER	Speech entity recognition
	KS	Keyword spotting
	IC	Intent classification
	SF	Slot filling
Sound	SAP	Speaker age prediction
	VSC	Vocal sound classification
	AAC	Automatic audio caption
	SEC	Sound event classification
	ASC	Acoustic scene classification
Music&Song	SED	Sound event detection with timestamps
	AQA	Audio question answering
	SID	Singer identification
	SMER	Singer and music emotion recognition
	MC	Music caption
Music&Song	MIC	Music instruments classification
	MNA	Music note analysis such as pitch, velocity
	MGR	Music genre recognition
	MD	Music distribution
	MSR	Music source separation

Qwen2.5-Audio.

Results.

- ! No Text Modality Metrics were reported.



Audio LLMs

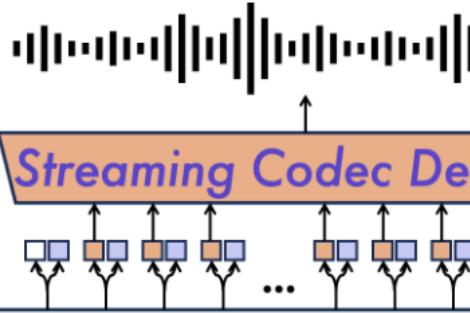
- Qwen2.5-Audio
- Qwen2.5-Omni ➡
- Moshi

Qwen2.5-Omni.

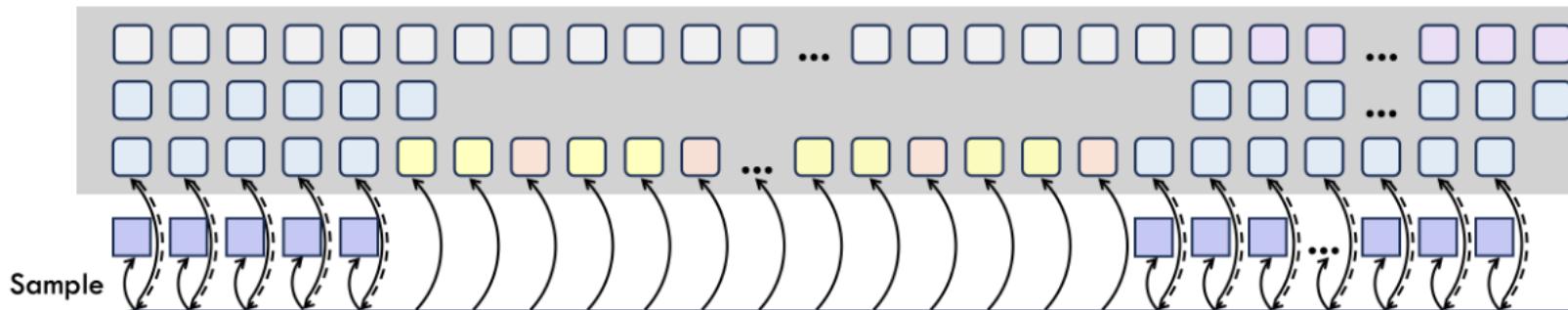
Architecture.

	Text Token
	Pad Token
	Codec Token
← Forward Propagation	
- → Backward Propagation	

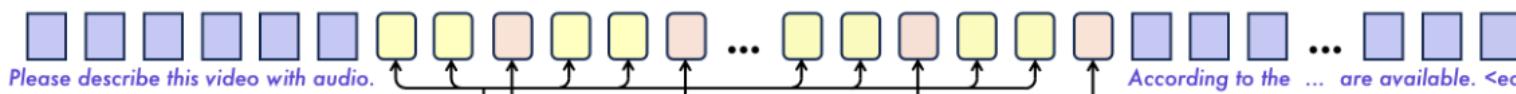
	Vision Hidden
	Text Hidden
	Audio Hidden
	Codec Hidden
	Pad Hidden



Qwen2.5-Omni Talker

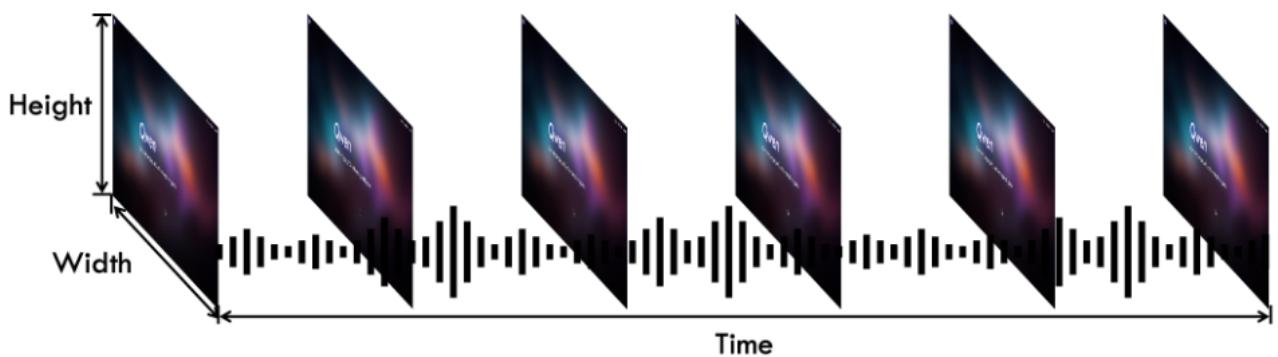


Qwen2.5-Omni Thinker



Vision Encoder

Audio Encoder



Qwen2.5-Omni.

Thinker Training.

3 training stages:

1. Encoders pretraining:



Audio Encoder.



Image Encoder.



LLM.

2. Finetuning:



Audio Encoder.



Image Encoder.



LLM.

3. Long Context Finetuning - 32k:



Audio Encoder.



Image Encoder.



LLM.

Qwen2.5-Omni.

Talker Training.

3 training stages:

1. Pretraining:

 Thinker.  Talker.

2. DPO:

 Thinker.  Talker.

3. SFT Instruction-tuning:

 Thinker.  Talker.

Qwen2.5-Omni.

Results.

Text Benchmarks.

Datasets	Qwen2-7B (text)	Qwen2-Audio	Qwen2.5-Omni.
MMLU*	69.3	33.2	65.6
CEval*	78.4	38.6	61.1
IFEval*	53.3	15.6	41.7
GSM8K*	82.3	18.4	85.4
Math23K*	92.3	23.0	87.1
Math401*	75.5	20.4	62.2

Qwen2.5-Omni.

Results.

Audio Understanding Benchmarks.

Table 8: Multimodality → Text performance of State-of-the-art and Qwen2.5-Omni

Datasets	Model	Performance
<i>Multimodal Understanding</i>		
OmniBench <i>Speech Sound Event Music Avg</i>	Gemini-1.5-Pro (Team et al., 2024)	42.67% 42.26% 46.23%
	MIO-Instruct (Wang et al., 2024g) (7B)	36.96% 33.58% 11.32%
	AnyGPT (7B) (Zhan et al., 2024)	17.77% 20.75% 13.21%
	video-SALMONN (13B) (Sun et al., 2024)	34.11% 31.70% 56.60%
	UnifiedIO2-xlarge (3.2B) (Lu et al., 2024a)	39.56% 36.98% 29.25%
	UnifiedIO2-xxlarge (6.8B) (Lu et al., 2024a)	34.24% 36.98% 24.53%
	MiniCPM-o (Yao et al., 2024)	- - - 40.5%
	Baichuan-Omni-1.5 (Li et al., 2025)	- - - 42.9%
	Qwen2.5-Omni-7B	55.25% 60.00% 52.83%

Qwen2.5-Omni.

Results.

Zero-shot Speech Generation Benchmarks.

Table 9: Zero-Shot Speech Generation

Datasets	Model	Performance		
<i>Content Consistency</i>				
SEED <i>test-zh test-en test-hard</i>	Seed-TTS _{ICL} (Anastassiou et al., 2024)	1.11	2.24	7.2
	Seed-TTS _{RL} (Anastassiou et al., 2024)	1.00	1.94	6.8
	MaskGCT (Wang et al., 2024e)	2.27	2.62	10.0
	E2 TTS (Eskimez et al., 2024)	1.97	2.19	-
	F5-TTS (Chen et al., 2024c)	1.56	1.83	8.3
	CosyVoice 2 (Du et al., 2024)	1.45	2.57	6.3
	CosyVoice 2-S (Du et al., 2024)	1.45	2.38	8.3
	Qwen2.5-Omni-7B _{ICL}	1.70	2.72	7.2
	Qwen2.5-Omni-7B _{RL}	1.42	2.33	6.3
<i>Speaker Similarity</i>				
SEED <i>test-zh test-en test-hard</i>	Seed-TTS _{ICL} (Anastassiou et al., 2024)	0.796	0.762	1.0
	Seed-TTS _{RL} (Anastassiou et al., 2024)	0.801	0.766	1.0
	MaskGCT (Wang et al., 2024e)	0.774	0.714	1.0
	E2 TTS (Eskimez et al., 2024)	0.730	0.710	1.0
	F5-TTS (Chen et al., 2024c)	0.741	0.647	1.0
	CosyVoice 2 (Du et al., 2024)	0.748	0.652	1.0
	CosyVoice 2-S (Du et al., 2024)	0.753	0.654	1.0
	Qwen2.5-Omni-7B _{ICL}	0.752	0.632	1.0
	Qwen2.5-Omni-7B _{RL}	0.754	0.641	1.0

Qwen2.5-Omni.

Questions.

- Qwen-TTS-Tokenizer (**issue**) ?
- Training Tasks and Data Details ?
- Why talker have to generate text tokens ? What about tasks interference ?

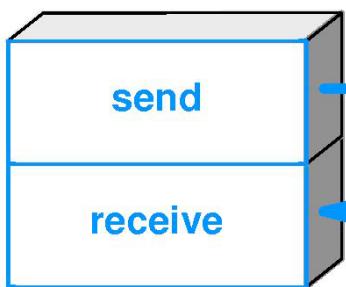
Audio LLMs

- Qwen2.5-Audio
- Qwen2.5-Omni
- Moshi ➡️

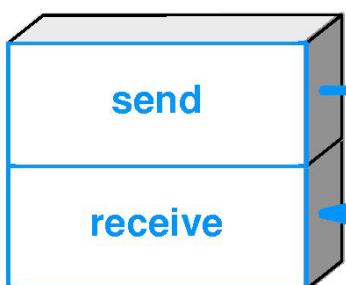
Simplex / Duplex / Half-duplex



(a) simplex



(b) full-duplex



(c) half-duplex

Moshi

Architecture.

- Pretrained LLM: Helium-7B
- Full-duplex

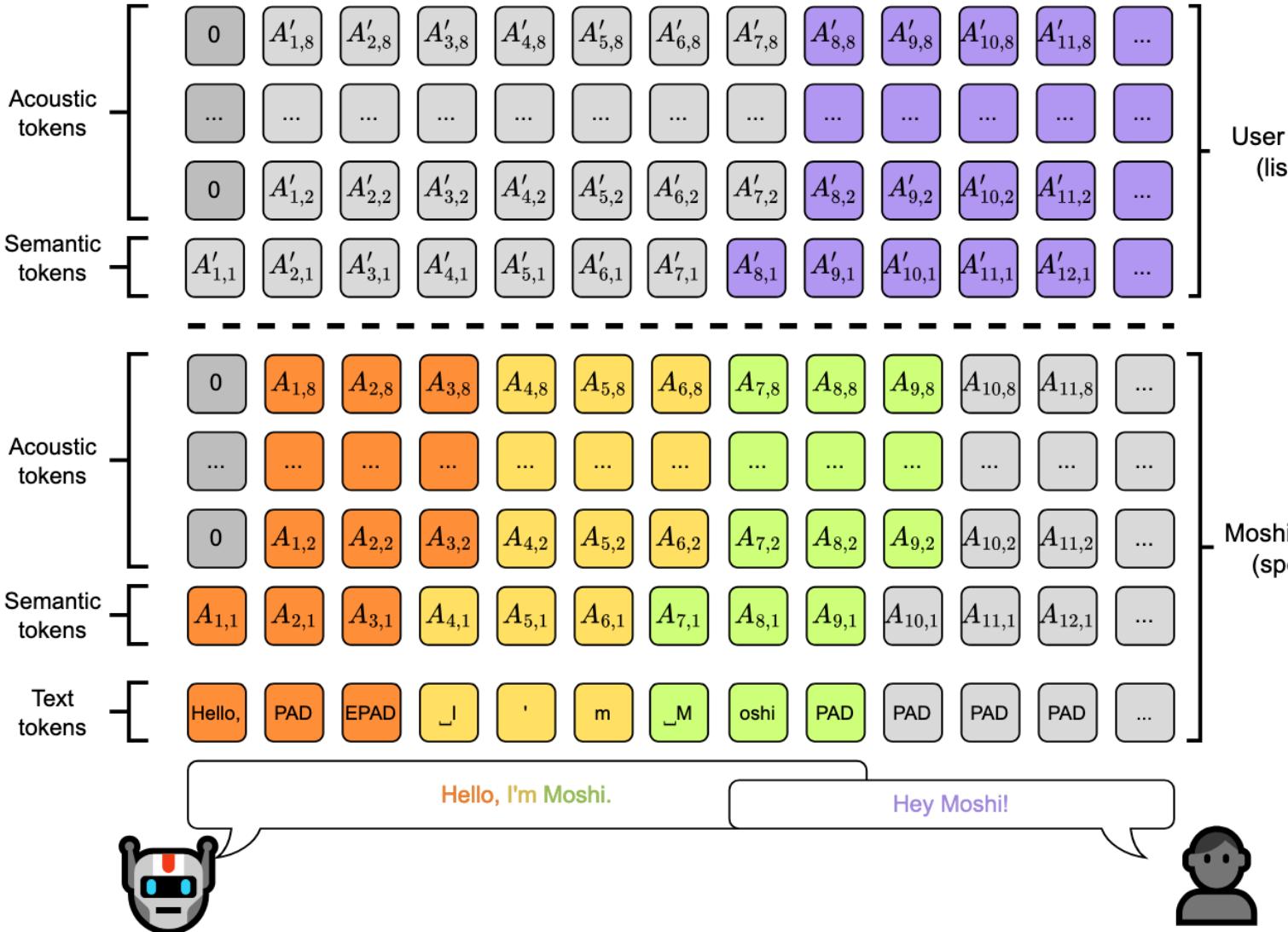


Figure 4: **Representation of the joint sequence modeled by Moshi.** Each row represents the tokens for a given step in the joint sequence ($V_{s,k}$) described in Eq. 1, with an acoustic delay $\tau = 1$, e.g. the input of the Temporal Transformer for the s -th token. Tokens are predicted from bottom to top in the Depth Transformer. At inference time, tokens under the dashed line (corresponding to Moshi) are sampled, while those above are fed from the user. This design allows for our model to handle overlapping speech.

Moshi

Training.

3 training stages:

1.  Audio Pretraining

Also text pretraining to prevent catastrophic forgetting

2.  Full-Duplex Training (Synth Data)

3.  Clean Dialogue Dataset SFT

Moshi Evaluation.

Table 7: **Performance of audio and text language modeling.** We report a based on scoring with negative log-likelihood, normalized by sequence length. L evaluated in a 5-shot setting. Reusing the terminology of Nguyen et al. (2024), \emptyset n unsupported modalities while - represents unreported numbers.

Model	Audio metrics			
	sWUGGY	sBLIMP	sTopic-StoryCloze	sStoryCloze
<i>Audio only - Cold Start</i>				
GSLM (Lakhotia et al., 2021)	64.8	54.2	66.6	53.3
AudioLM (Borsos et al., 2022)	71.5	64.7	-	-
TWIST (Hassid et al., 2023)	72.2	56.5	-	-
Moshi	74.8	59.9	80.9	56.9
<i>Audio only - Warm Start</i>				
TWIST (Hassid et al., 2023)	74.5	59.2	76.4	55.4
VoxtLM (Maiti et al., 2023)	62.9	53.9	-	-
Spirit-LM (Nguyen et al., 2024)	69.5	58.0	72.9	54.8
Moshi	74.3	58.9	81.8	58.7
<i>Text and audio - Warm Start</i>				
VoxtLM (Maiti et al., 2023)	66.1	57.1	-	-
Spirit-LM (Nguyen et al., 2024)	69.0	58.3	82.9	61.0
Moshi after single-stream pretraining	72.6	58.8	83.0	60.8
Moshi after multi-stream instruct	63.0	55.2	83.6	62.7
Moshi after multi-stream instruct, synthetic voice	60.9	54.6	82.5	60.9

Comparison

	Qwen2.5-Audio	Qwen2.5-Omni	Moshi (O)
Audio-In	✓	✓	✓
Audio-Out		✓	✓
Image-In		✓	
Full-Duplex			✓

Materials

- Moshi overview SpeechInfo [1/2], [2/2]
- **Sesame Conversational Voice**