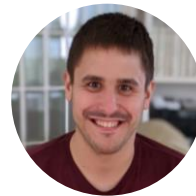


CLIPsonic: Text-to-Audio Synthesis with Unlabeled Videos and Pretrained Language-Vision Models

Hao-Wen Dong^{1,2*} Xiaoyu Liu¹ Jordi Pons¹ Gautam Bhattacharya¹
Santiago Pascual¹ Joan Serrà¹ Taylor Berg-Kirkpatrick² Julian McAuley²

¹ Dolby Laboratories ² University of California San Diego

* Work done during an internship at Dolby



Overview – Text-to-Audio Synthesis

(These samples are generated by our proposed model.)

More samples



salu133445.github.io/clipsonic

Prior Work – Text-to-Audio Synthesis

- Diffsound (Yang et al., 2023)
- AudioGen (Kreuk et al., 2023)
- AudioLDM (Liu et al., 2023)
- Make-An-Audio (Huang et al., 2023)
- Noise2Music (Huang et al., 2023)
- MusicLM (Agostinelli et al., 2023)



All rely on large amounts of
text-audio training pairs

Can we learn text-to-audio synthesis
without using any text-audio pairs?

Yang et al., "Diffsound: Discrete Diffusion Model for Text-to-sound Generation," *TASLP*, 2022.

Kreuk et al., "AudioGen: Textually Guided Audio Generation," *ICLR*, 2023.

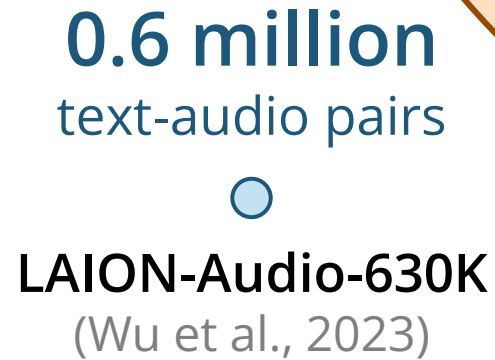
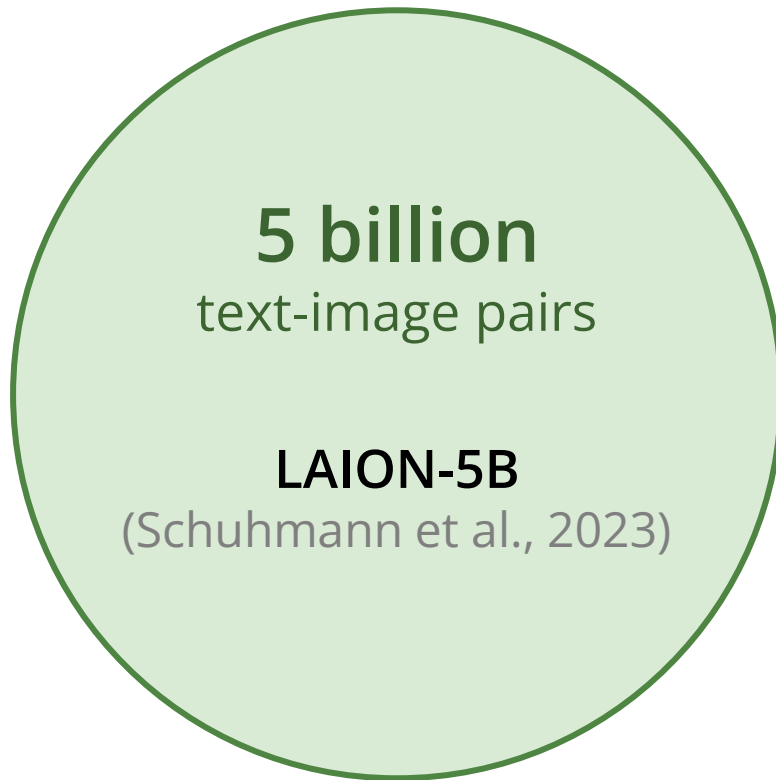
Liu et al., "AudioLDM: Text-to-Audio Generation with Latent Diffusion Models," *ICML*, 2023.

Huang et al., "Make-An-Audio: Text-To-Audio Generation with Prompt-Enhanced Diffusion Models," *ICML*, 2023.

Huang et al., "Noise2Music: Text-conditioned Music Generation with Diffusion Models," *arXiv preprint arXiv:2302.03917*, 2023.

Agostinelli et al., "MusicLM: Generating Music From Text," *arXiv preprint arXiv:2302.03917*, 2023.

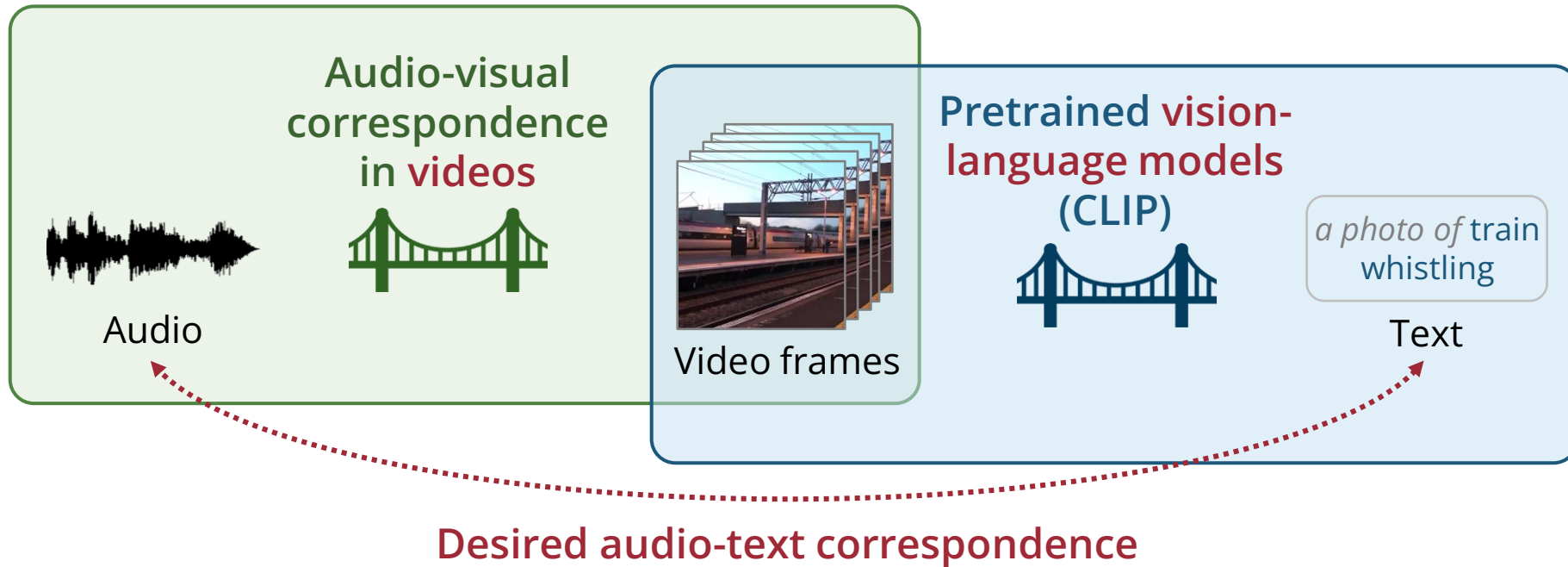
Why NOT Text-audio Pairs?



YouTube videos!

500 hours of videos
uploaded per minute

Leveraging the Visual Domain as a Bridge



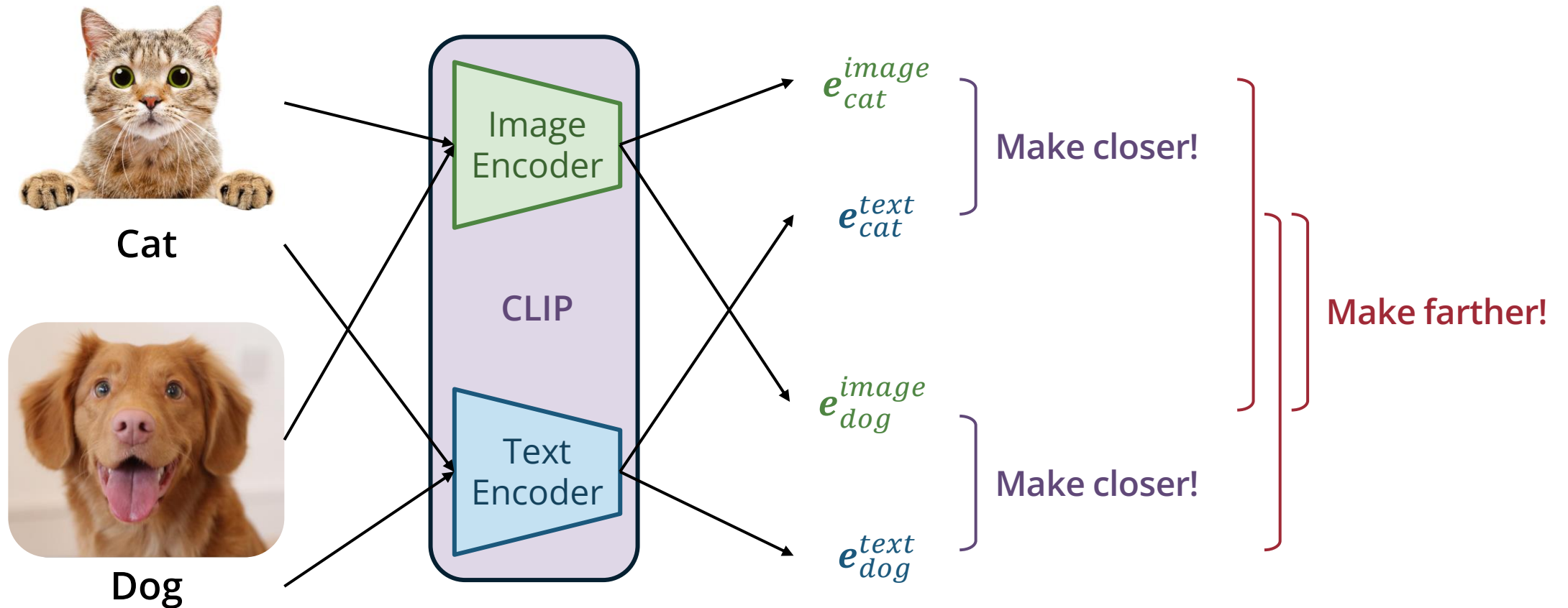
No text-audio pairs required!

Scalable to large video datasets!

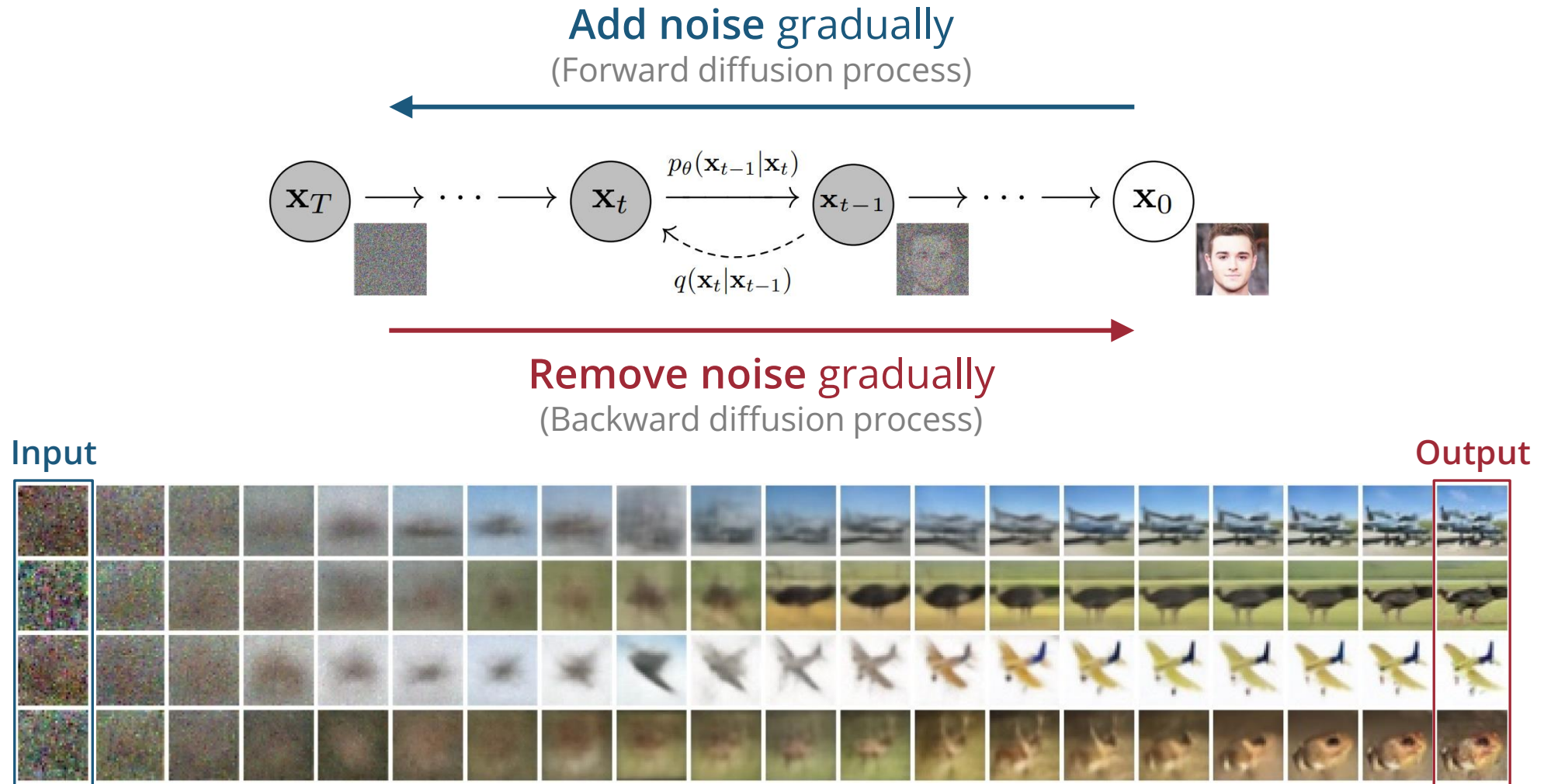
Background

CLIP (Contrastive Language-Image Pretraining)

- Learn a **shared embedding space** for images and texts via *contrastive learning*



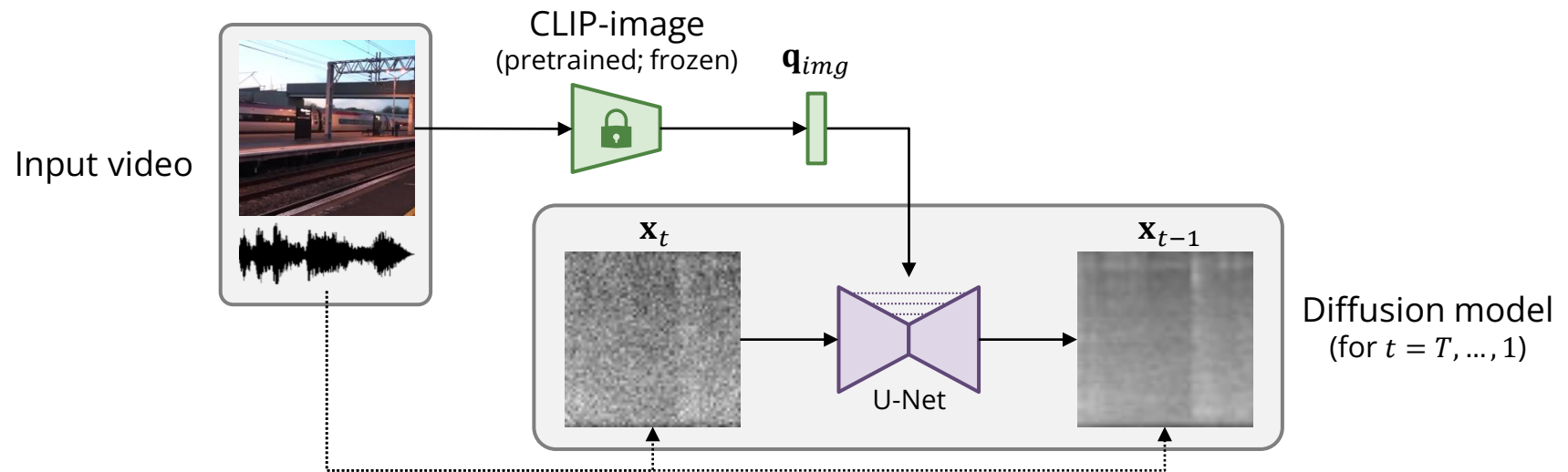
Diffusion Model



Method

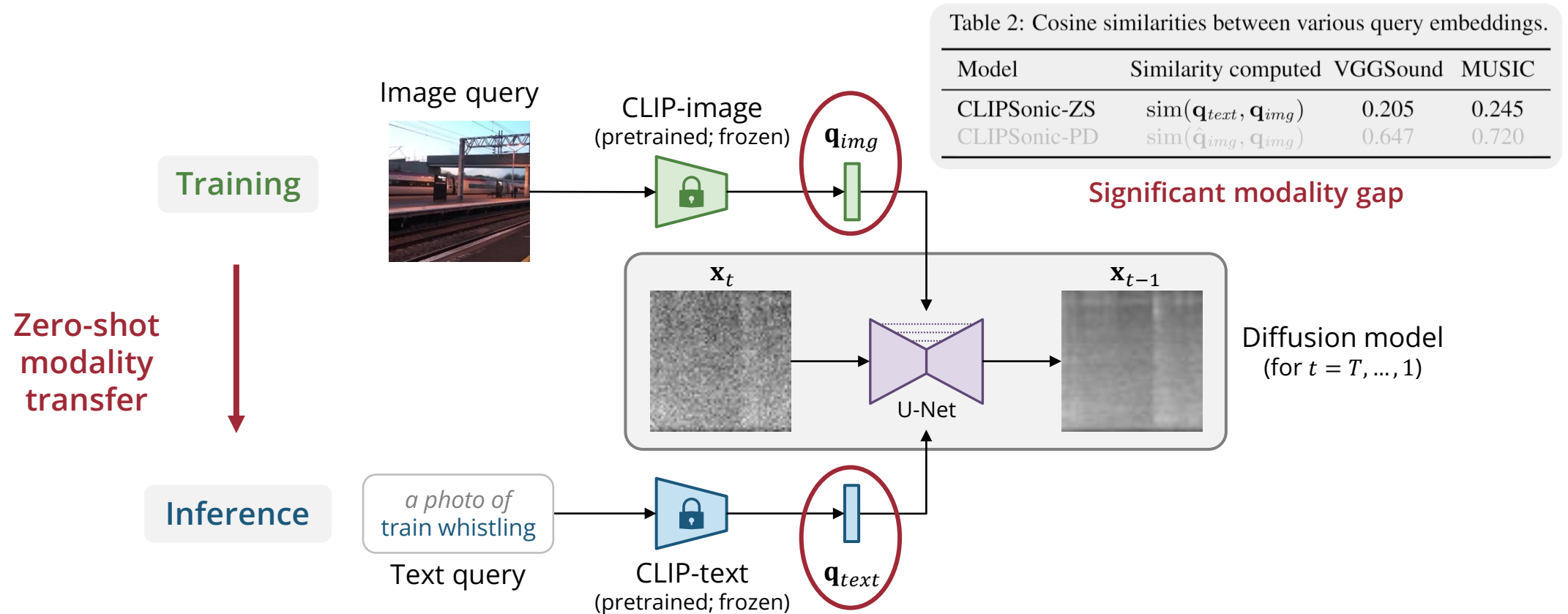
Training – Image-queried

- We train an image-to-audio synthesis model using a diffusion model on mel spectrograms and a pretrained CLIP-image encoder



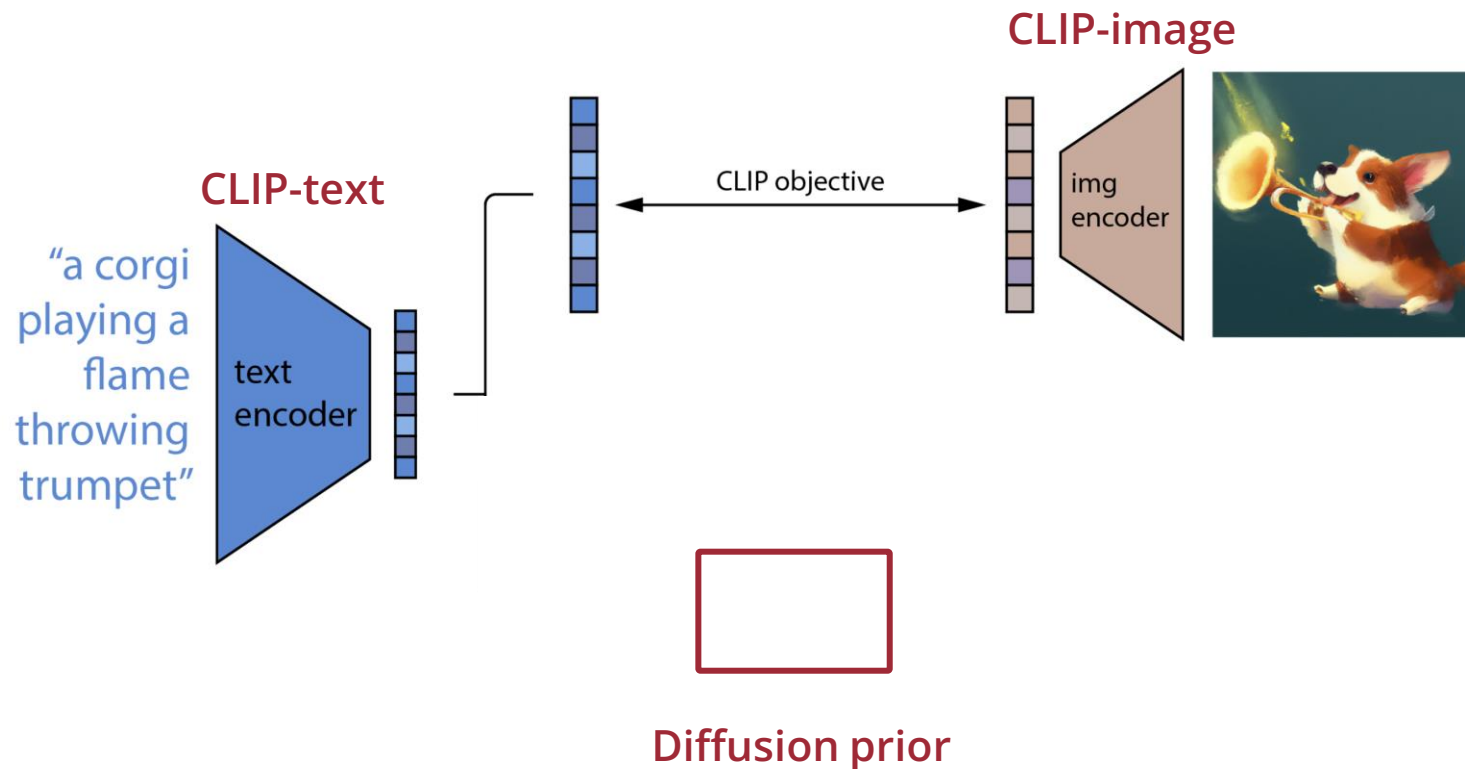
Inference – Zero-shot Modality Transfer (CLIP Sonic-ZS)

- We first explore using a pretrained CLIP-text encoder directly



How to overcome this modality gap?

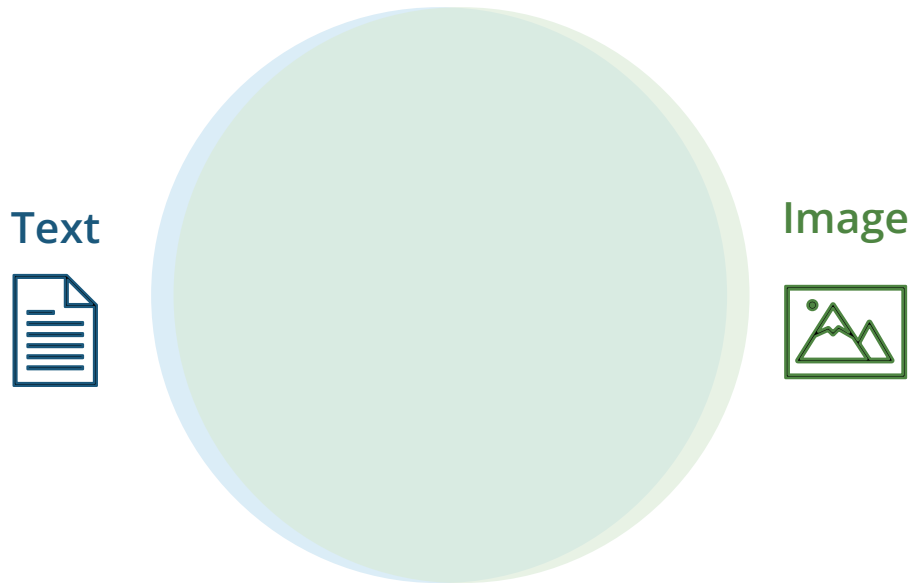
- We leverage a pretrained diffusion prior model (Ramesh et al., 2022)



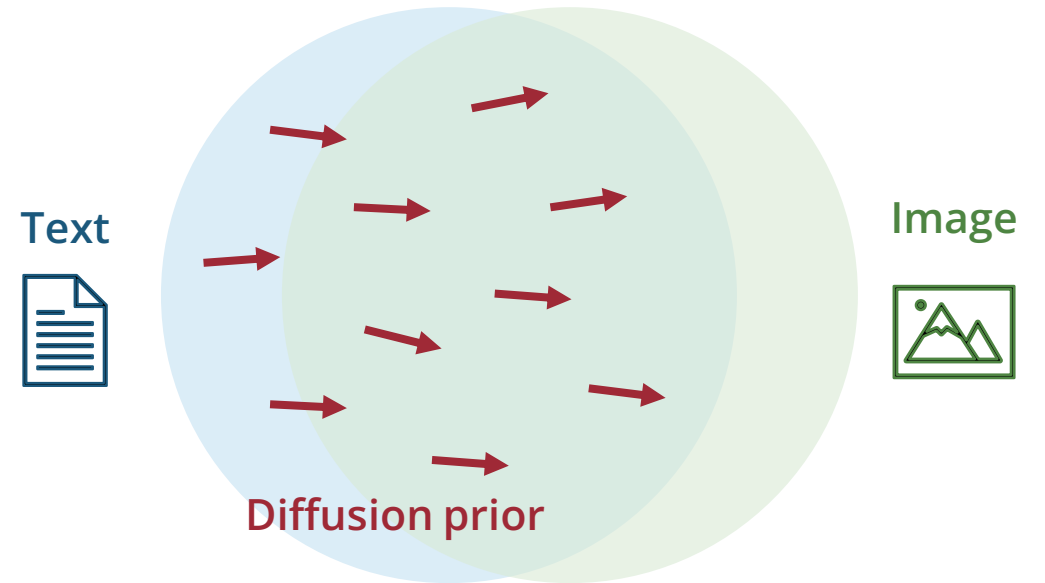
Diffusion Prior (Ramesh et al., 2022)

CLIP embedding spaces

Ideal case

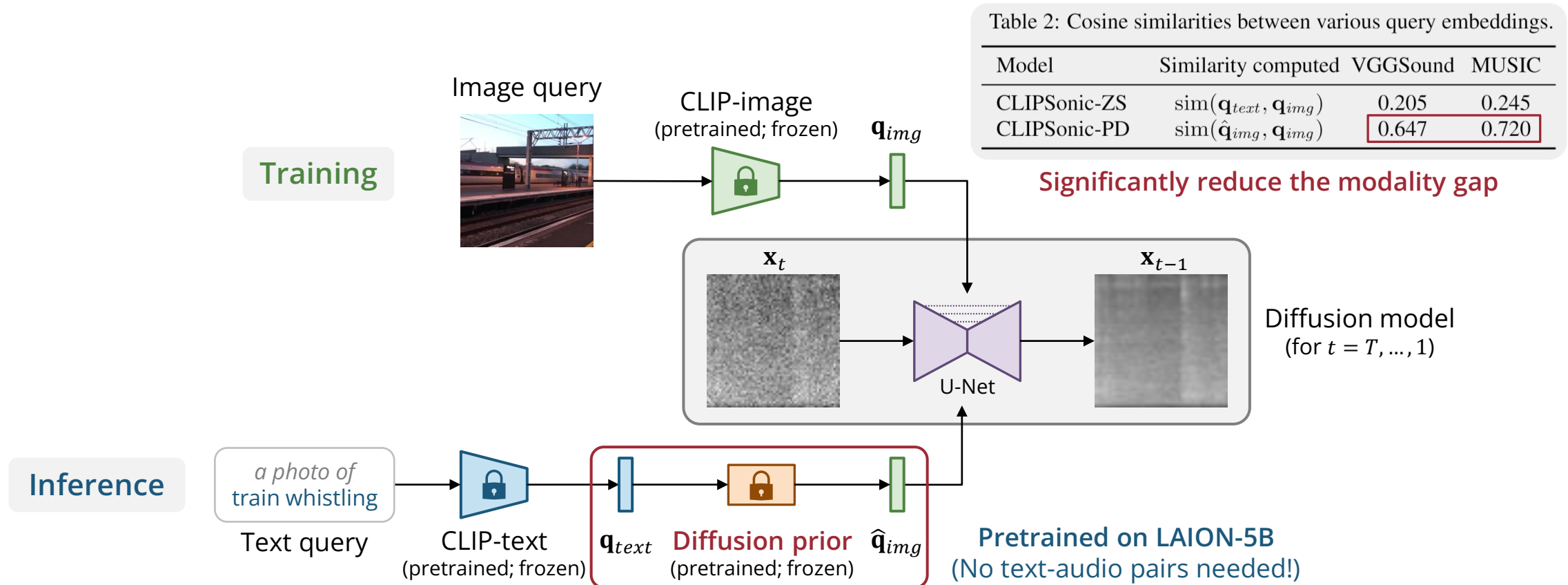


In practice



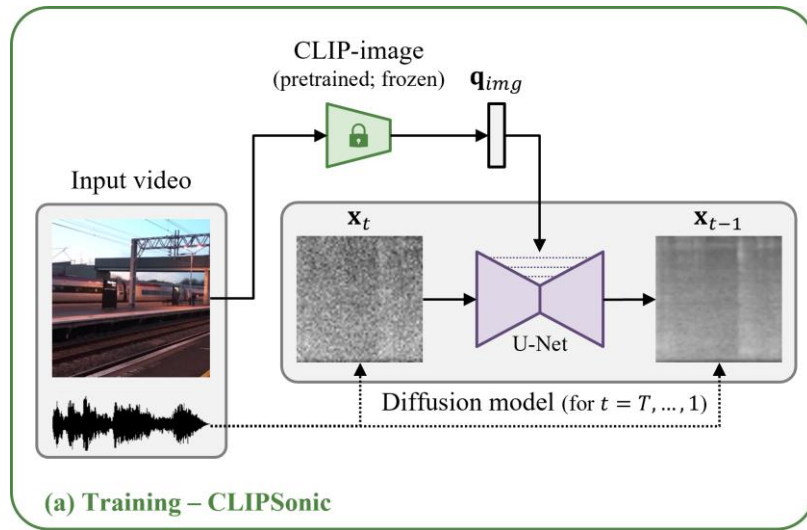
Inference – Pretrained Diffusion Prior (CLIP Sonic-PD)

- We then explore using a **pretrained diffusion prior model** (Ramesh et al., 2022)



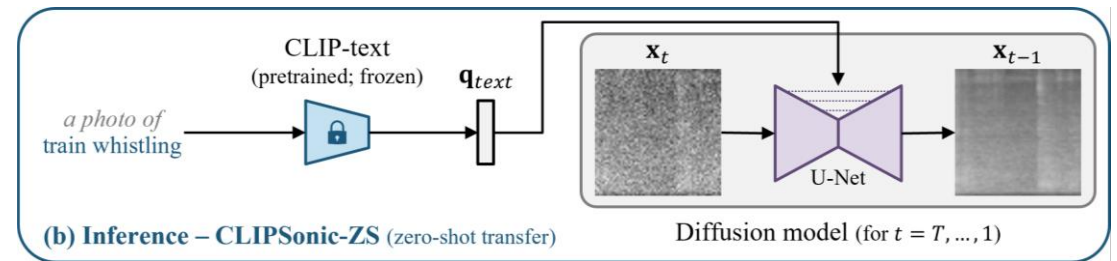
Recap

Training

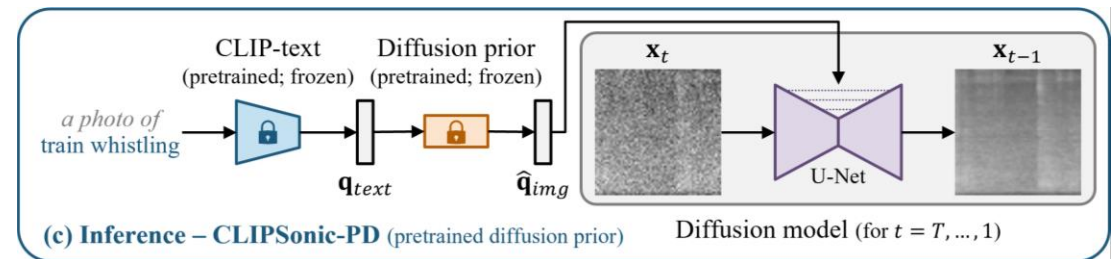


CLIP Sonic-IQ
(Image-queried)

Inference



CLIP Sonic-ZS
(Zero-shot transfer)



CLIP Sonic-PD
(Pretrained diffusion prior)

Experiments

Data

MUSIC

(Zhao et al., 2018)



Violin



Acoustic guitar



Accordion

Music instrument playing videos

VGGSound

(Chen et al., 2020)



Hedge trimmer running



Dog bow-wow



Bird chirping, tweeting

Noisy videos with diverse sounds

Implementation Details

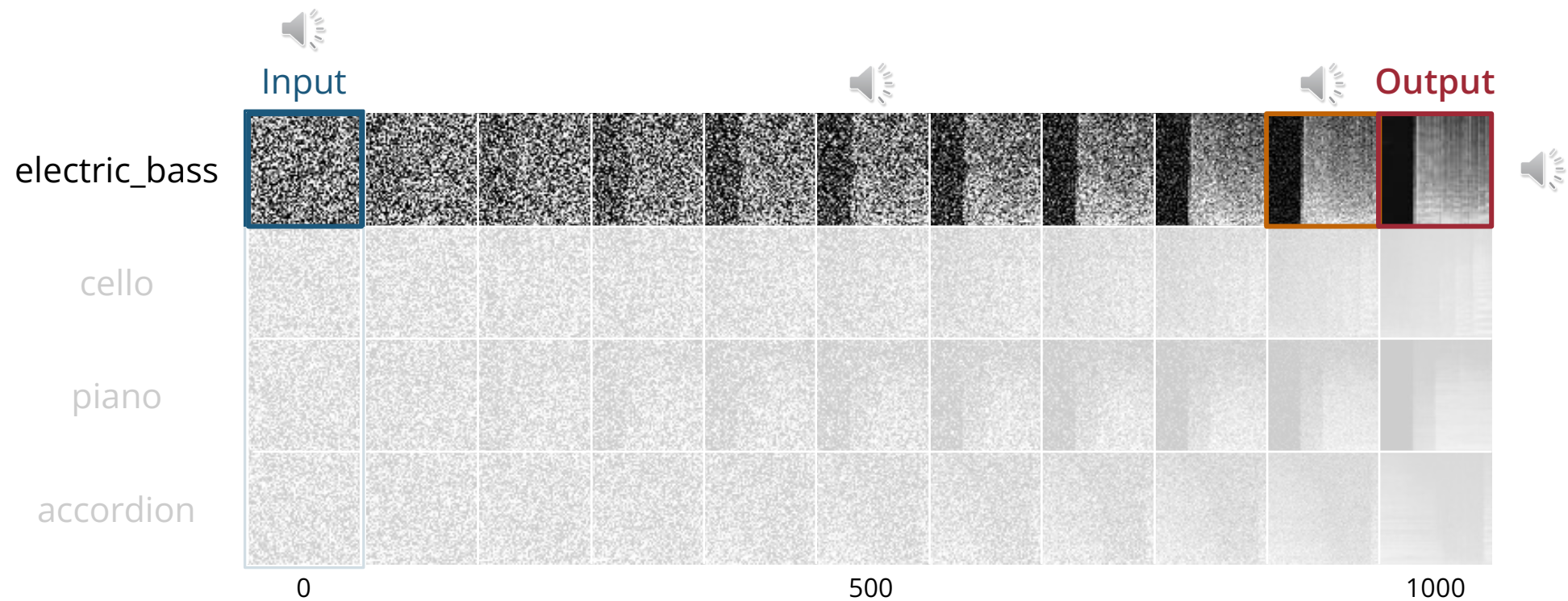
Mel spectrogram configuration

- Sampling rate: 16 kHz
- Hop size: 512
- FFT filter size: 2048
- 64 mel bands
- Inverted back to waveforms using **BigVGAN** (Lee et al., 2023)

Diffusion model

- Based on **Improved DDPM** (Nichol and Dhariwal, 2019)
- Diffusion steps:
 - Training: 4000
 - Inference: 1000
- Training iterations
 - MUSIC: 200K (1 day on 2 RTX 2080 Tis)
 - VGGSound: 500K (2 days on 2 RTX 2080 Tis)

Inference – Examples



Text-to-Audio Synthesis – Demo

Rapping



Sea waves



Thunder



Smoke detector beeping



Playing table tennis



Playing violin fiddle



Text-to-Audio Synthesis – Demo

Rapping

Sea waves

Smoke detector beeping

CLIPsonic-ZS

(zero-shot modality transfer)



CLIPsonic-PD

(pretrained diffusion prior)



Playing table tennis

Thunder

Playing violin fiddle

CLIPsonic-ZS

(zero-shot modality transfer)



CLIPsonic-PD

(pretrained diffusion prior)



Text-to-Audio Synthesis – Listening Test

Table 3: Listening test results for text-to-audio synthesis (MOS).

Model	VGGSound		MUSIC	
	Fidelity	Relevance	Fidelity	Relevance
CLIPSonic-ZS	2.55 ± 0.22	2.01 ± 0.27	2.98 ± 0.23	3.87 ± 0.24
CLIPSonic-PD	3.04 ± 0.20	2.86 ± 0.25	3.67 ± 0.18	3.91 ± 0.24
Ground truth	3.78 ± 0.19	3.54 ± 0.29	3.90 ± 0.17	4.34 ± 0.18

Significant performance improvement against the baseline!

Image-to-Audio Synthesis – Demo (Out-of-distribution)

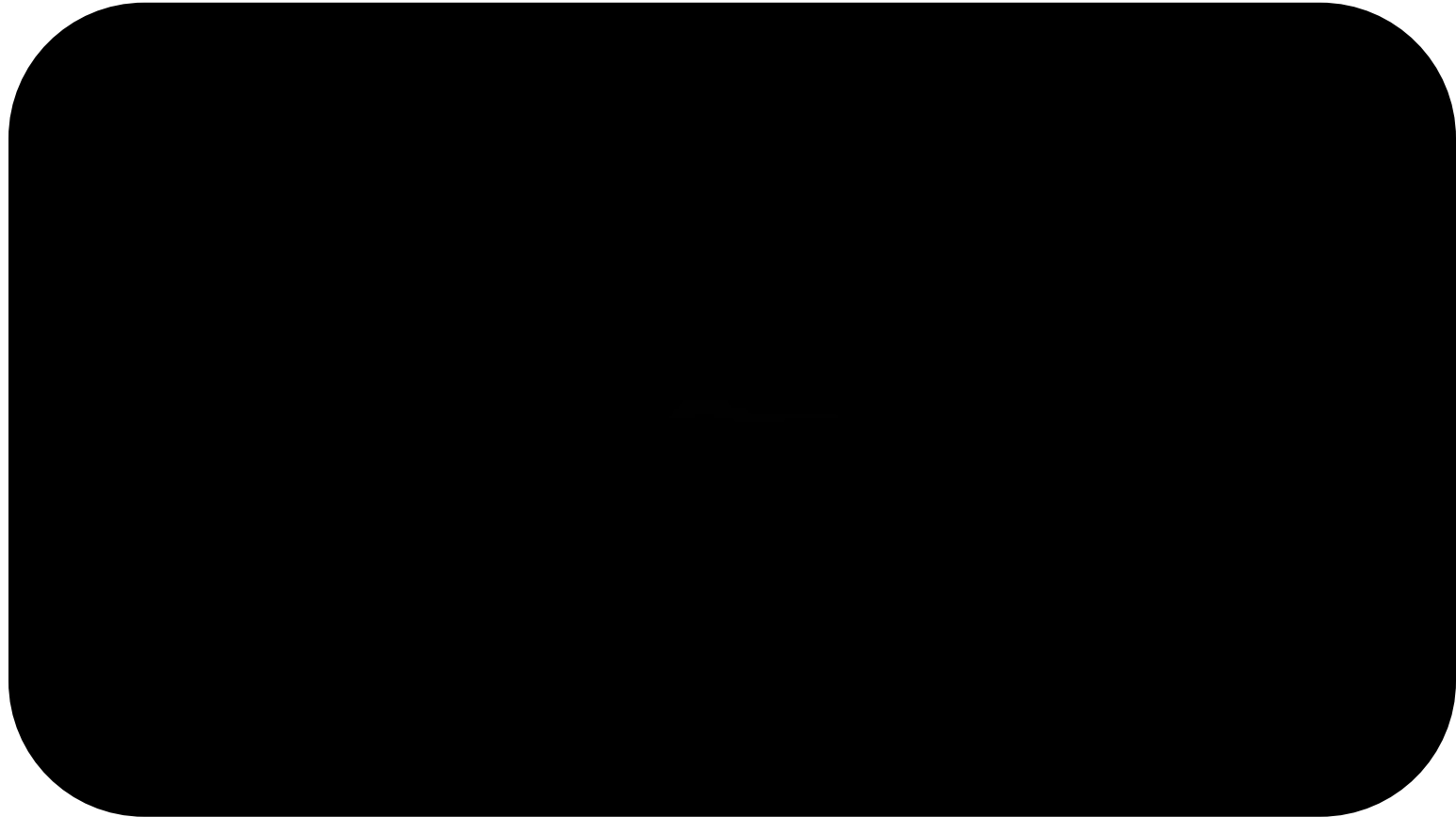


Image-to-Audio Synthesis – Demo (Out-of-distribution)



CLIPSONIC-IQ (ours)



Im2wav (Sheffer & Adi, 2023)



SpecVQGAN (Iashin & Rahtu, 2021)



Image-to-Audio Synthesis – Listening Test

Table 4: Listening test results for image-to-audio synthesis (MOS).

Model	Fidelity	Relevance
CLIPSonic-IQ (image-queried)	3.29 ± 0.16	3.80 ± 0.19
SpecVQGAN [20]	2.15 ± 0.17	2.54 ± 0.23
im2wav [21]	2.19 ± 0.15	3.90 ± 0.22

State-of-the-art image-to-audio performance!

Objective Evaluation Metrics

- Evaluated with Fréchet audio distance (FAD) and CLAP score

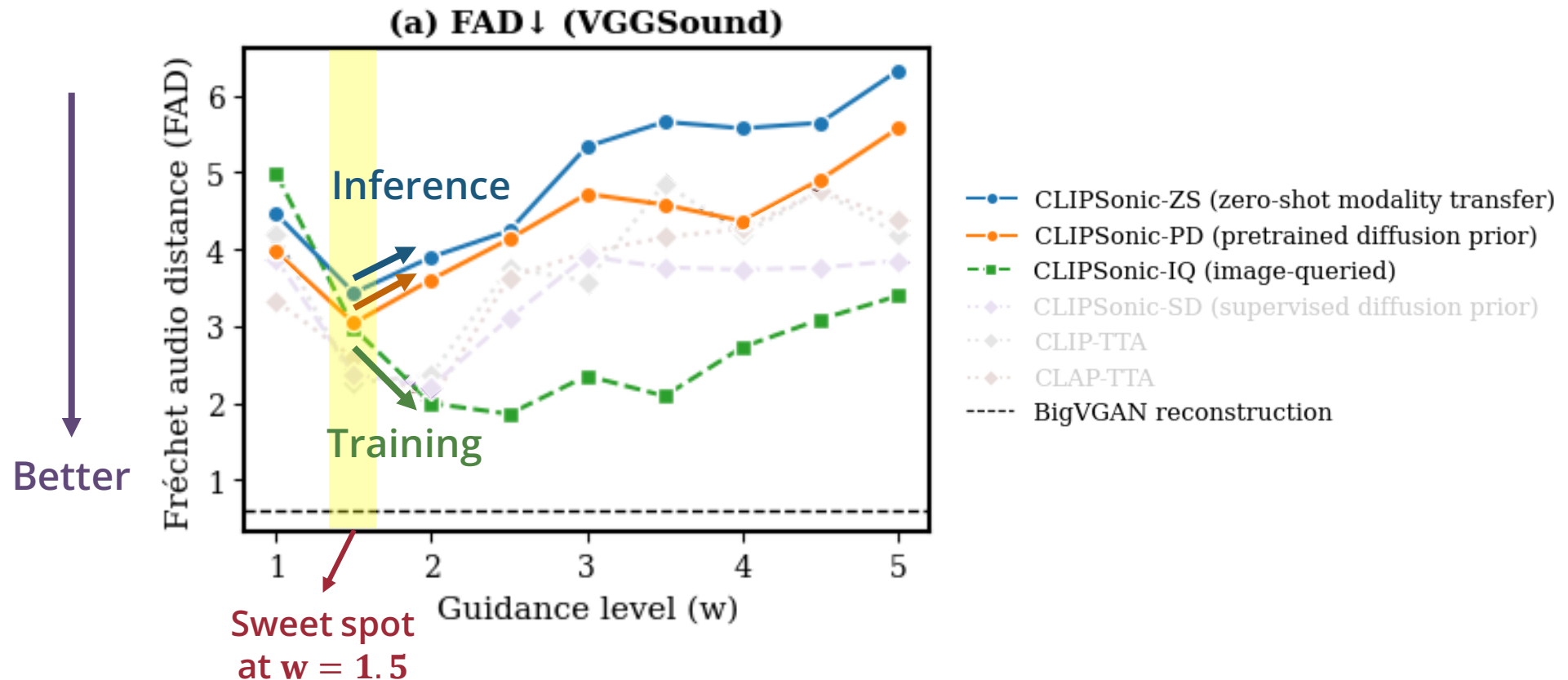
Table 1: Evaluation results on VGGSound and MUSIC datasets, evaluated at $w = 1.5$.

Model	Without text-audio pairs	Query modality		VGGSound		MUSIC	
		Training	Inference	FAD ↓	CLAP score ↑		
CLIPSONIC-IQ (image-queried)	-	Image	Image	2.37	0.234	12.13	0.299
CLIPSONIC-ZS (zero-shot modality transfer)	✓	-	-	2.26	0.292	9.39	0.298
CLIPSONIC-PD (pretrained diffusion model)	✓	-	-	2.58	0.296	10.92	0.303
Big vGAN mel spectrogram reconstruction	-	-	-	0.60	0.204	6.21	0.272

Check out our paper for more results!

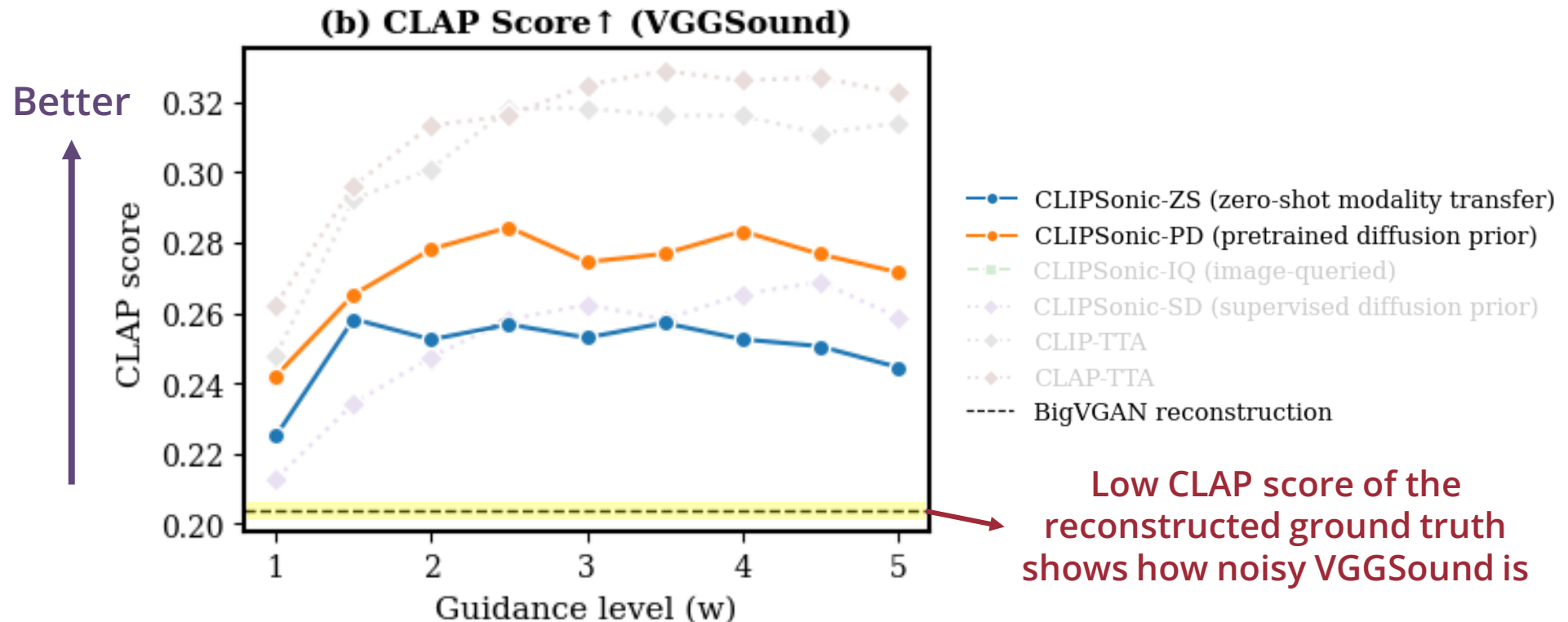
Effects of Classifier-free Guidance

- A guidance level of $w = 1.5$ leads to the lowest FADs for inference



Effects of Classifier-free Guidance

- Larger guidance level leads to stronger adherence to input text query

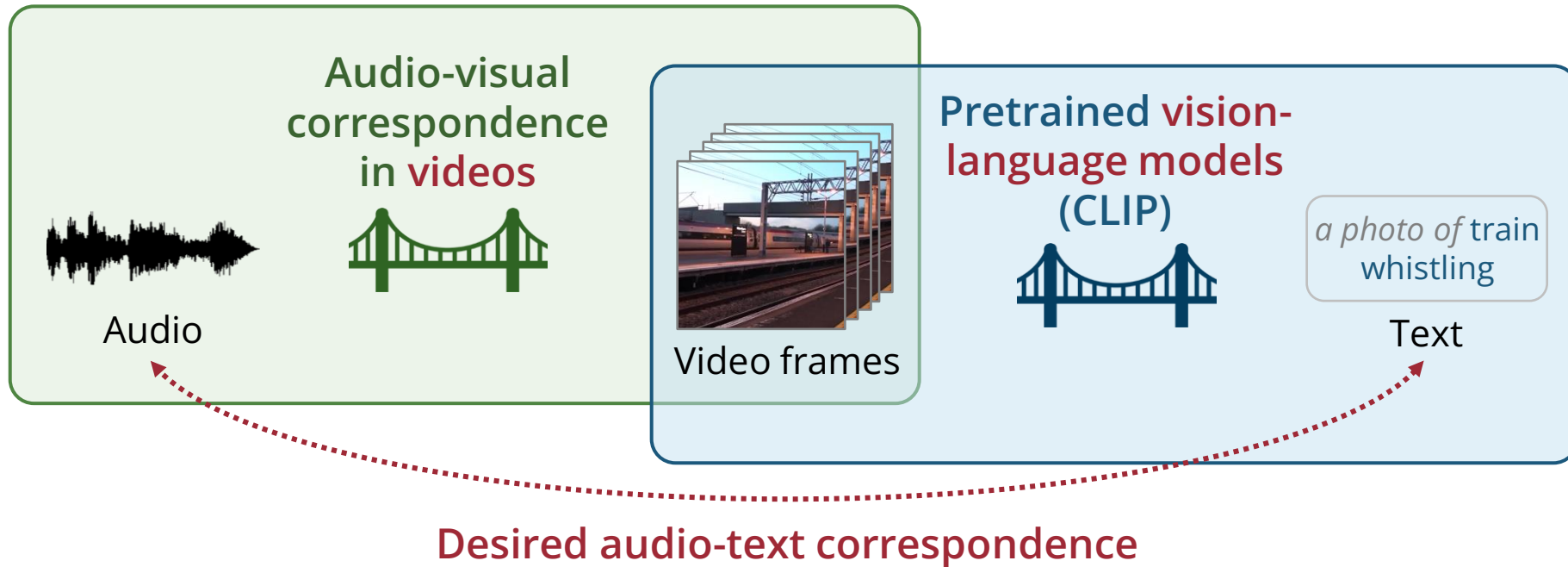


Limitations & Future Work

- Off-screen sounds occur frequently in videos
- Cannot handle purely audio-specific queries
- Can we enable compositional prompts?
- Scale up to larger video datasets!

Conclusion

Leveraging the Visual Domain as a Bridge

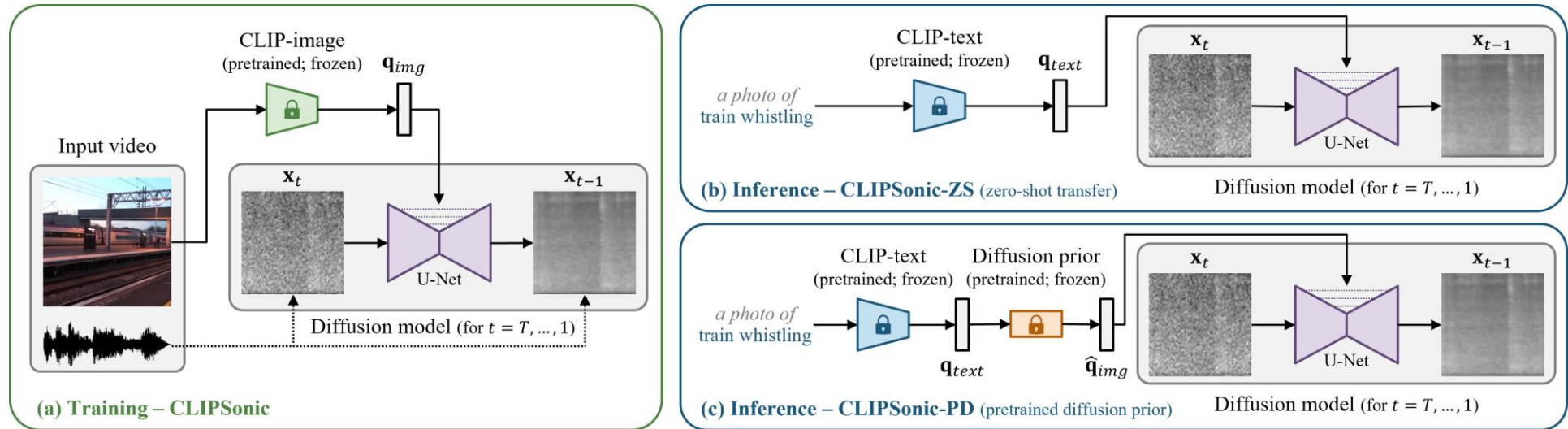


No text-audio pairs required!

Scalable to large video datasets!

Summary

- Proposed a text-to-audio synthesis model that **requires no text-audio pairs**
- Achieved strong text-to-audio synthesis performance
- Achieved state-of-the-art performance in image-to-audio synthesis



Thank you!

Paper: arxiv.org/abs/2306.09635
Demo: salu133445.github.io/clipsonic

