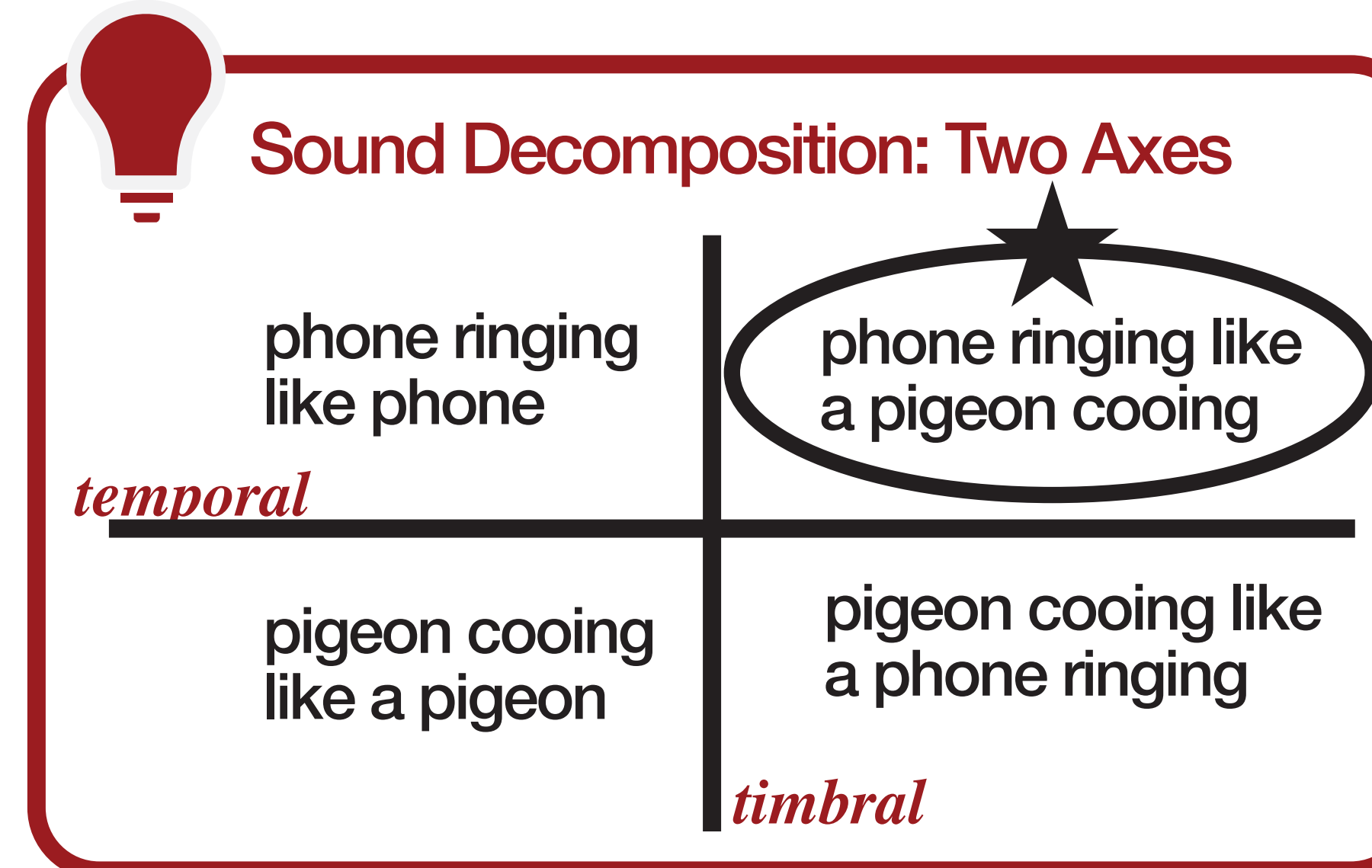


# MIX2MORPH: LEARNING SOUND MORPHING FROM NOISY MIXES

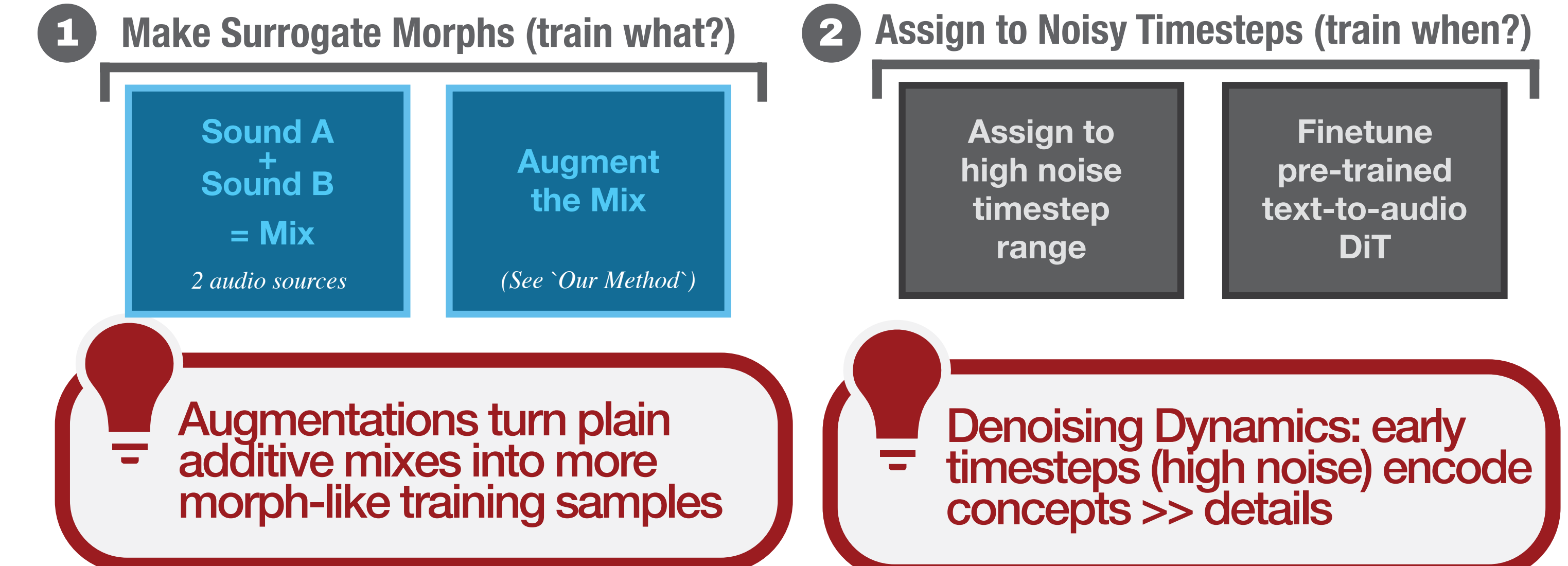
Annie Chu<sup>1\*2</sup>, Hugo Flores García<sup>1,2</sup>, Oriol Nieto<sup>1</sup>, Justin Salamon<sup>1</sup>, Bryan Pardo<sup>2</sup>, Prem Seetharaman<sup>1</sup>



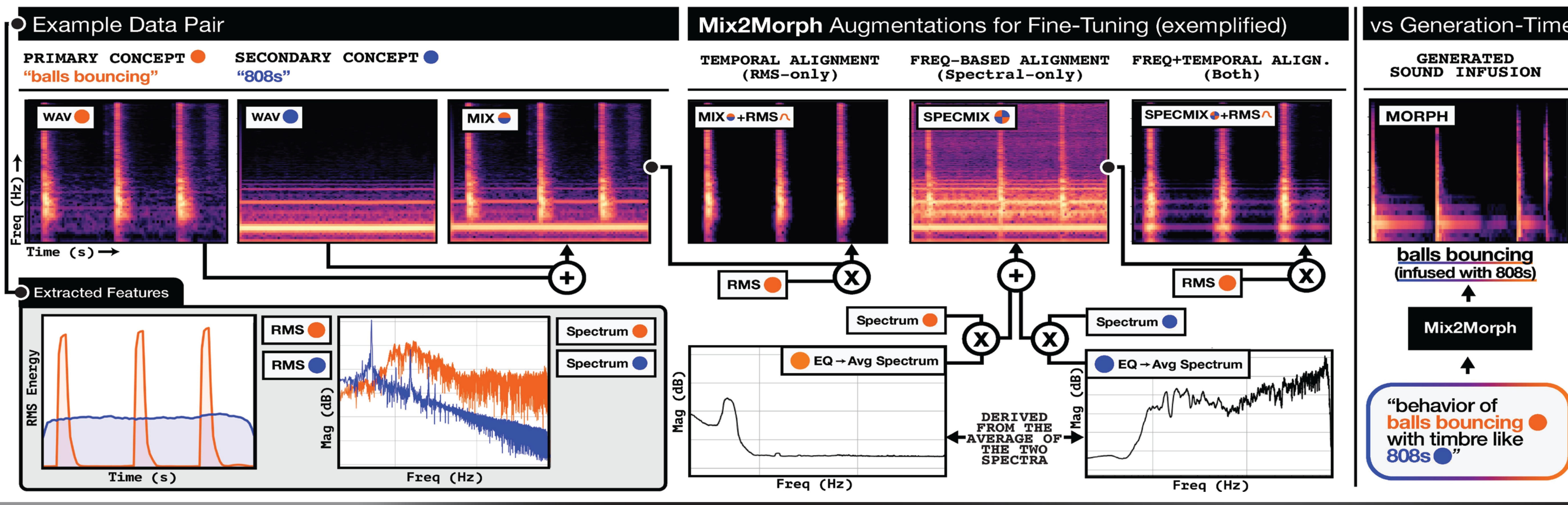
## GOAL: Asymmetric Sound Morphs (aka Sound Infusions)



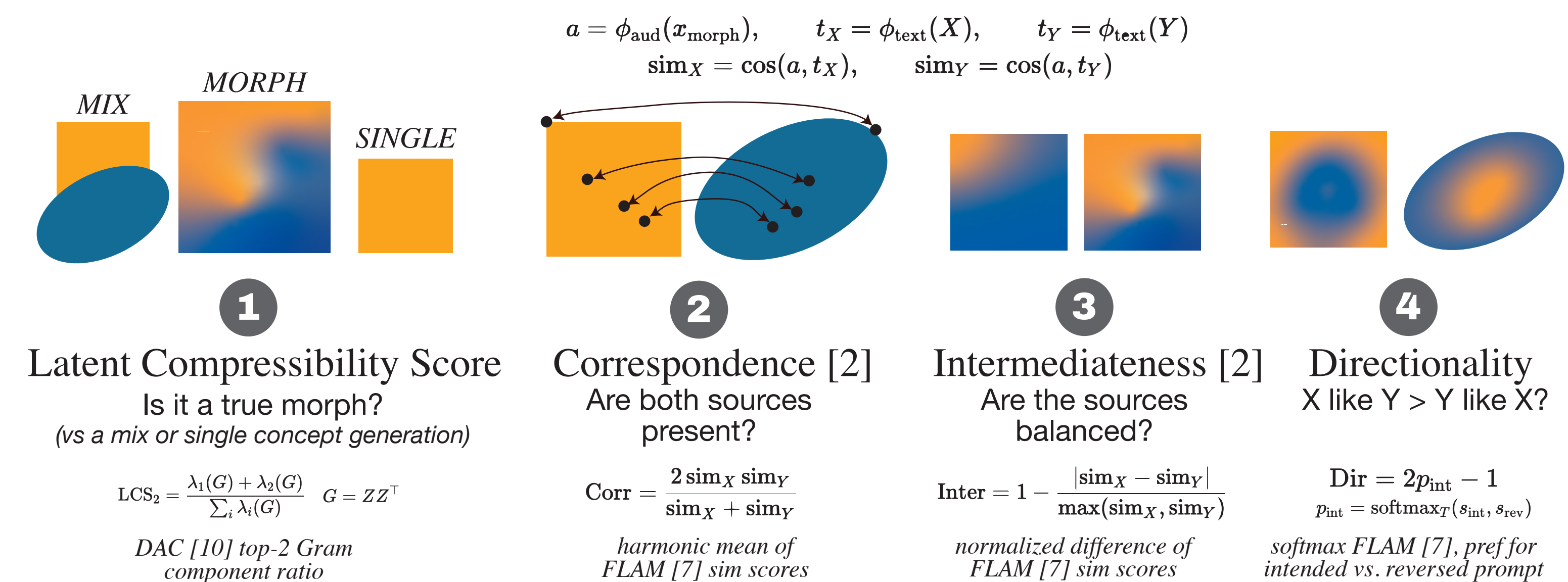
## How? Finetune DiT on Noisy Mixes



## Our Method: Augment Mixes to Approximate Morphs, Assign @ High Ts



## Evaluating Sound Infusions: Morph Metrics



## References

[1] D. Maul, "From Smart Phones to Carrier Pigeons (when more tech is the last thing you need!)," *Derek Maul: Words & Photographs for the Journey*, Jan. 19, 2022. [Online]. Available: <https://derekmaul.blog/2022/01/19/from-smart-phones-to-carrier-pigeons-when-more-tech-is-the-last-thing-you-need/>.

[2] M. F. Caetano and N. Osaka, "A formal evaluation framework for sound morphing," in *Non-Cochlear Sound: Proc. 38th Int. Computer Music Conf. (ICMC)*, Ljubljana, Slovenia, Sept. 9-14, 2012. Michigan Publishing, 2012.

[3] L. M. Heller and L. Wolf, "When hybrid sound effects are better than real recordings," *Proc. Meetings Acoust.*, vol. 46, no. 1, Art. no. 050002, 2022, doi: 10.1121/2.0001581.

[4] A. J. Oxenham, "How we hear: The perception and neural coding of sound," *Annu. Rev. Psychol.*, vol. 69, no. 1, pp. 27-50, Jan. 2018, doi: 10.1146/annurev-psych-122216-011635.

[5] Z. M. Smith, B. Delgutte, and A. J. Oxenham, "Chimaeric sounds reveal dichotomies in auditory perception," *Nature*, vol. 416, no. 6876, pp. 87-90, Mar. 2002, doi: 10.1038/416087a.

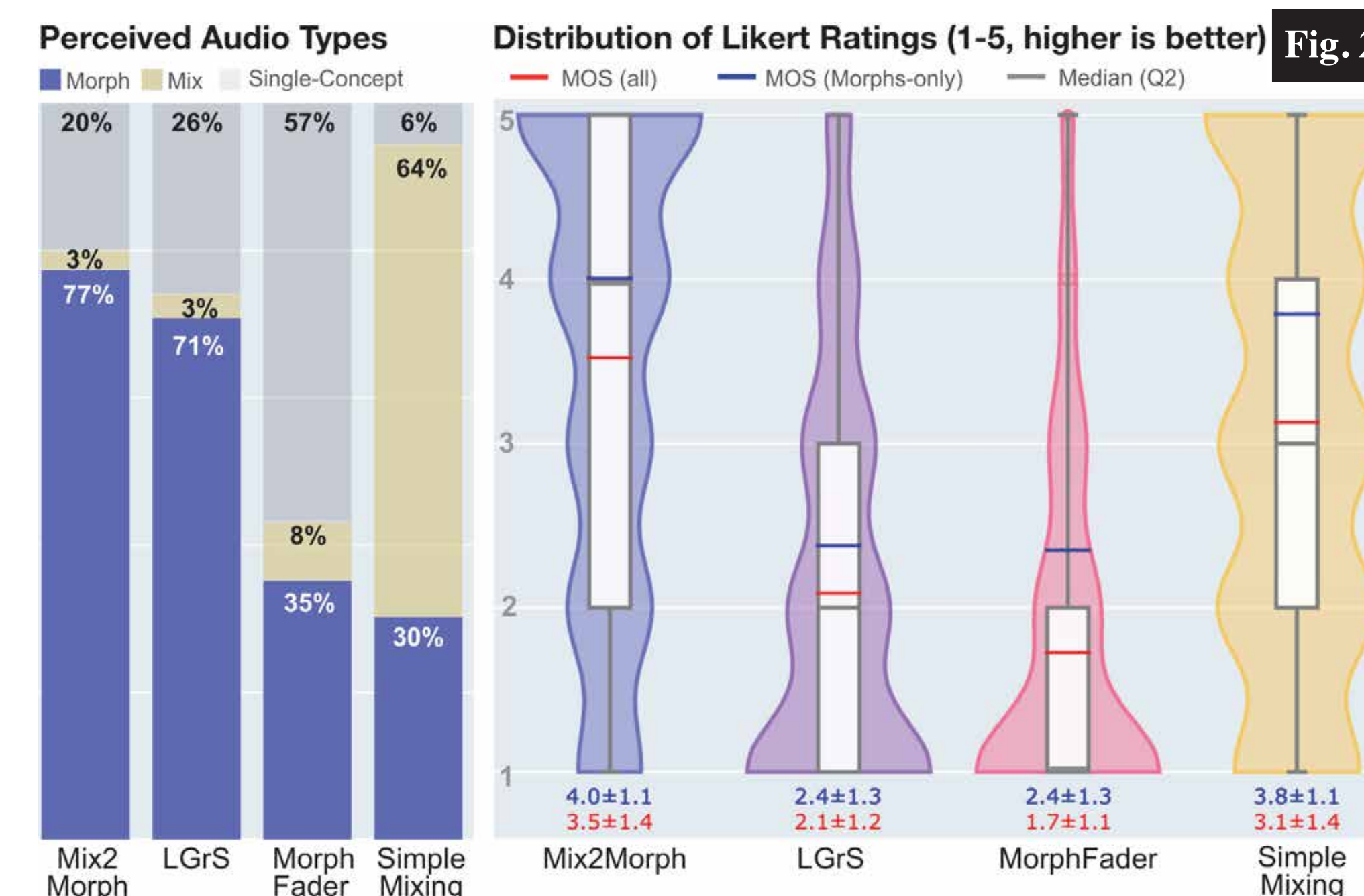
[6] P. Kamath, C. Gupta, and S. Nanayakkara, "MorphFader: Enabling fine-grained controllable morphing with text-to-audio models," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, 2025, pp. 1-5, doi: 10.1109/ICASSP49660.2025.10890164.

[7] Y. Wu, C. Tsirigotis, K. Chen, C.-Z. A. Huang, A. Courville, O. Nieto, P. Seetharaman, and J. Salamon, "FLAM: Frame-wise language-audio modeling," in *Proc. 42nd Int. Conf. Mach. Learn. (ICML)*, Proc. Mach. Learn. Res., vol. 267, pp. 67719-67740, 2025.

[8] N. Tokui and T. Baker, "Latent granular resynthesis using neural audio codecs," in *Extended Abstracts for the Late-Breaking Demo Session of the 26th Int. Society for Music Information Retrieval Conf. (ISMIR)*, Daejeon, South Korea, 2025, doi: 10.48550/arXiv.2507.19202.

[9] X. Niu, J. Zhang, and C. P. Martin, "SoundMorpher: Perceptually-uniform sound morphing with diffusion model," arXiv:2410.02144, 2024, doi: 10.48550/arXiv.2410.02144.

[10] Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar, "High-fidelity audio compression with improved rvqgan," *Advances in Neural Information Processing Systems*, vol. 36, pp. 27980-27993, 2023.



## Mix2Morph ablations

Best Mix2Morph Model 3-way augmentation at [0.5, 1]

## vs Baselines

Mix2Morph >> More perceptually coherent midpoint sound infusions

Table 1

Model	t_start	t_end	RMS	Spectral	Both	None	LCS ↑	Corr. ↑	Inter. ↑	Direct. ↑	FAD ↓
base	-	-	-	-	-	-	0.136	0.678	0.611	0.525	1.219
+Timestep Allocation	0	1	✓	✗	✗	✗	0.128	0.699	0.646	0.173	1.230
	0.25	1	✓	✗	✗	✗	0.143	0.705	0.658	0.278	1.235
	0.5	1	✓	✗	✗	✗	0.141	0.721	<b>0.672</b>	0.296	1.221
	0.75	1	✓	✗	✗	✗	0.134	0.717	0.653	0.364	1.225
+Augmentation Mode	0.5	1	✓	✗	✓	✗	0.135	0.700	0.623	0.363	1.226
	0.5	1	✓	✓	✓	✗	<b>0.150</b>	<b>0.725</b>	0.648	<b>0.436</b>	<b>1.220</b>
	0.5	1	✓	✓	✓	✓	0.143	0.712	0.650	0.349	1.222
Simple Mixing	-	-	-	-	-	-	0.132	<b>0.758</b>	<b>0.690</b>	-9.25e-13	1.293
LGrS	-	-	-	-	-	-	0.173	0.539	0.638	-0.119	1.290
MorphFader	-	-	-	-	-	-	0.085	0.418	0.421	-9.72e-13	1.430
SoundMorpher	-	-	-	-	-	-	<b>0.242</b>	0.591	0.641	-9.64e-13	1.380
Mix2Morph	0.5	1	✓	✓	✓	✗	0.150	0.725	0.648	<b>0.436</b>	<b>1.220</b>

## Results: Listening Study (N=25) and Objective Metrics