



Speech Recovery for Real-World Self-Powered Intermittent Devices

Poster #: 3750

Yu-Chen Lin^{1,3}, Tsun-An Hsieh³, Kuo-Hsuan Hung³, Cheng Yu³, Harinath Garudadri⁵, Yu Tsao³, Tei-Wei Kuo^{1,2,4}

¹Department of Computer Science and Information Engineering, National Taiwan University, Taipei, Taiwan

²NTU High Performance and Scientific Computing Center, National Taiwan University, Taipei, Taiwan

³Research Center for Information Technology Innovation, Academia Sinica, Taipei, Taiwan

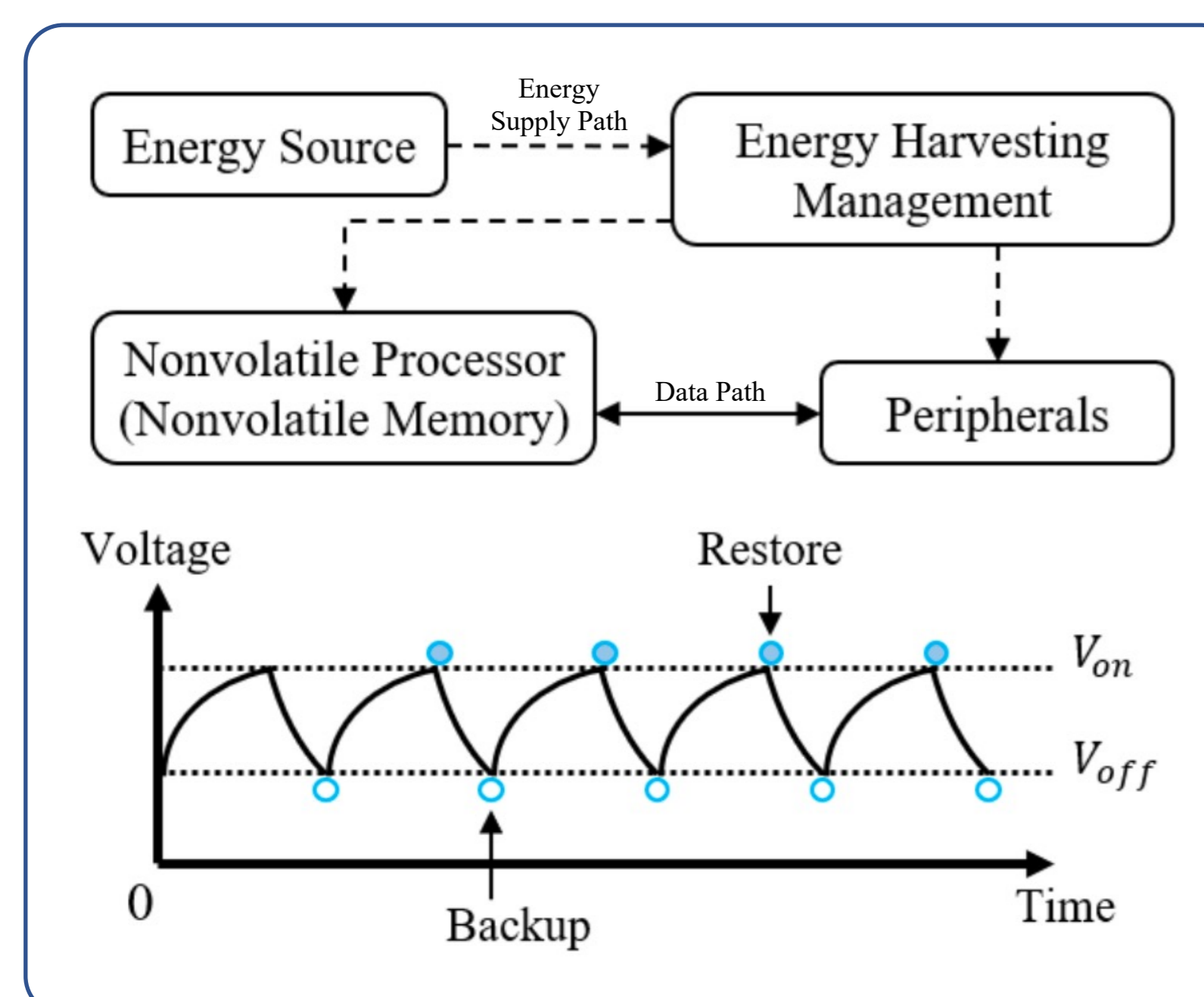
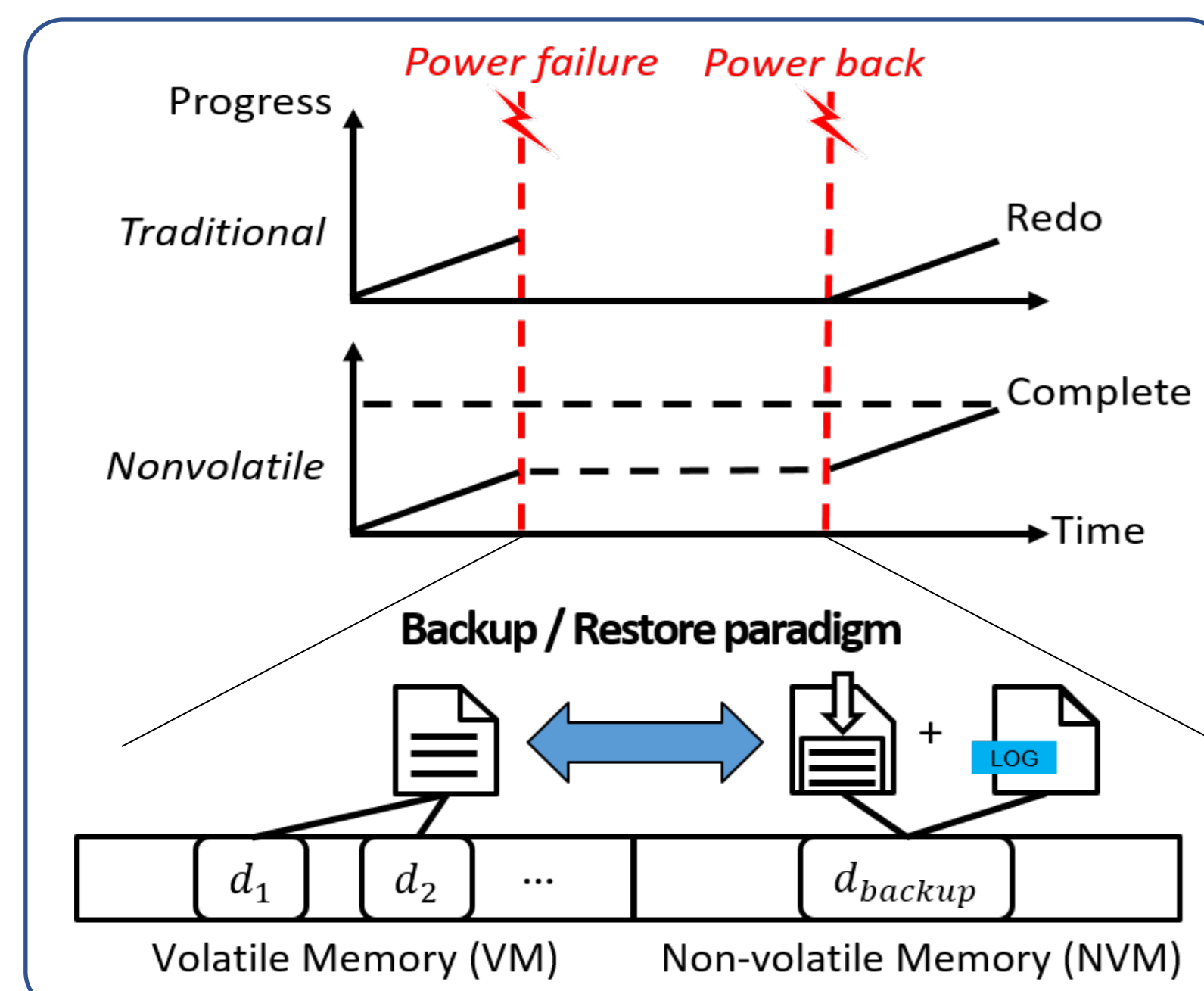
⁴Department of Computer Science, City University of Hong Kong, Hong Kong

⁵Qualcomm Institute, University of California, San Diego, USA

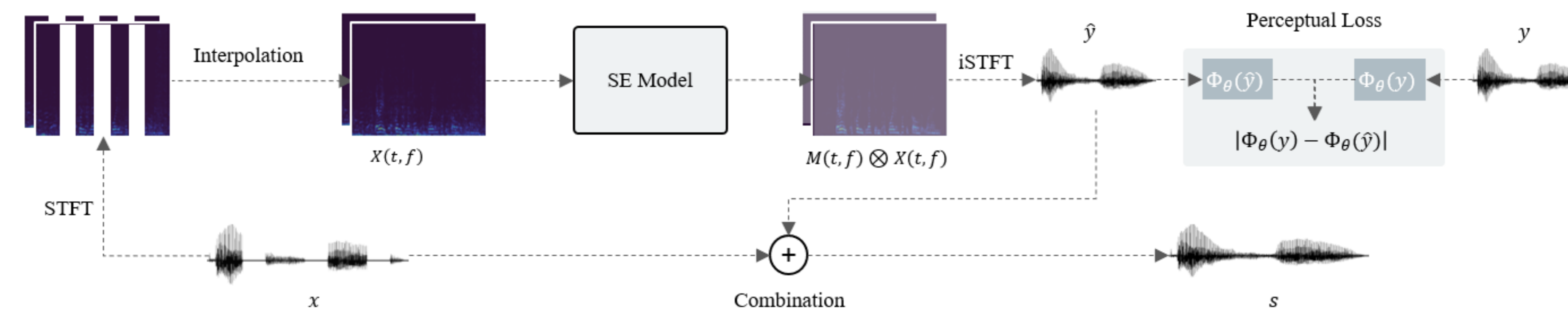


Introduction

- The Era of IoT Devices
 - Wearable devices outnumbers the worlds popularity
 - High cost in maintenance batteries in IoT devices — recharge/pollution
- Alternatives: Energy Harvesting Devices
 - Volatile data will be lost frequently due to power failures
 - Systems need to be frequently recovered after power resumption
- Self-Powered Intermittent System
 - Redoing in battery-powered systems
 - Intermittent completion by *preserving forward progress* at runtime



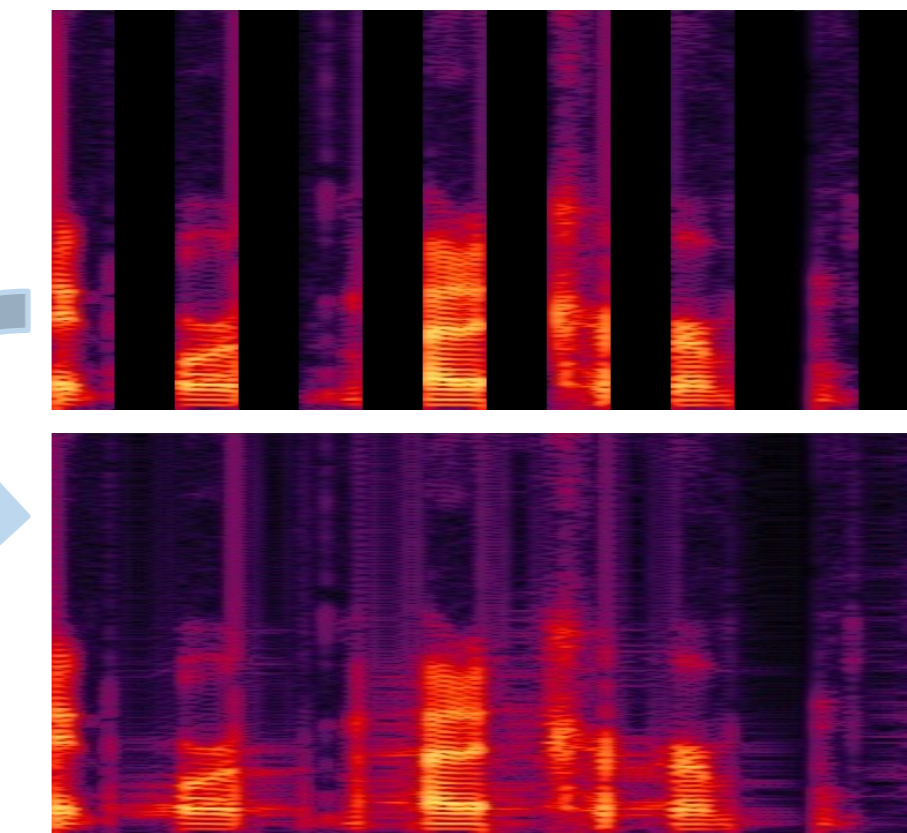
Intermittent Speech Recovery



- Null Segment Interpolation
 - Frequency domain interpolation
 - Frequency components are more observable while being interpolated in the frequency domain
 - Interpolation policy
 - Weighting ratio

$$r(t) = \frac{t - (t_1 - 1)}{(t_2 + 1) - (t_1 - 1)}$$

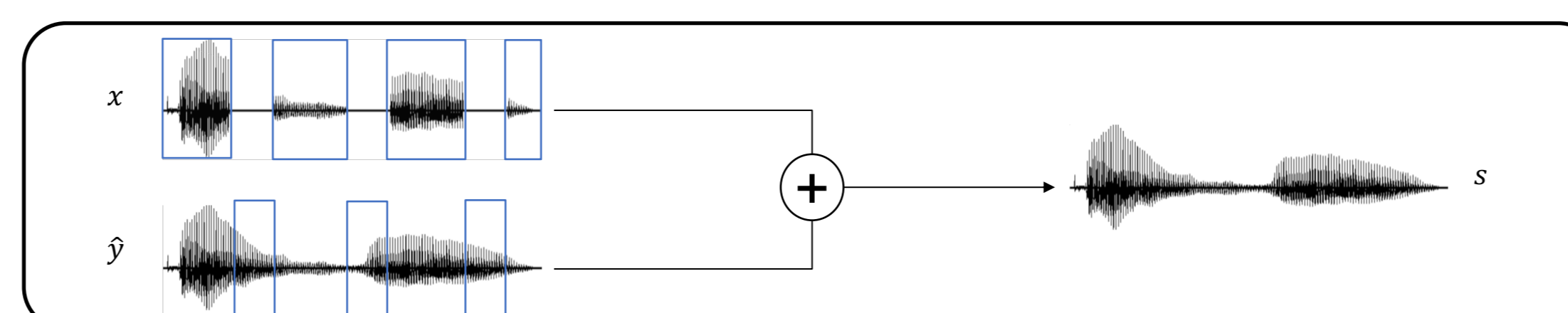
$$X(t, f) = (1 - r(t))X(t_1 - 1, f) + r(t)X(t_2 + 1, f)$$



- Interpolated Speech Refinement via Speech Enhancement
 - Deep complex U-Net architecture
 - Phase information preservation
 - Perceptual/feature losses are proven being able to
 - Improve speech quality and intelligibility
 - Render phonetic and speaker information

$$\mathcal{L}(y, \hat{y}) = \frac{1}{CL} \sum_{c=1}^C \sum_{l=1}^L (|\Phi_{\theta}(y)_{c,l} - \Phi_{\theta}(\hat{y})_{c,l}|^p)^{1/p}$$

- Combination
 - Enhancement model changes values in non-null segments
 - Combination policy
 - \hat{y} is the estimated speech generated by the complex U-Net
 - x is the original intermittent speech



Experimental Results

- Intermittent Recording System Setup

Intermittent Microphone Device	
Restore threshold (V_{on})	2.8 V
Backup threshold (V_{off})	2.3 V
Recording energy consumption	5.6 mW
EHM & Power Supply	
Capacitance	200 μ F
Training energy source	1.5 to 5.5 mW (with a step of 0.25)
Testing energy source	2.0 to 5.0 mW (with a step of 1.0)

- Performance Evaluation

Energy source	Period (ms)			Intermittent			Interpolated			ISR+MSE			ISR+PL		
	On	Off		PESQ	STOI	WER	PESQ	STOI	WER	PESQ	STOI	WER	PESQ	STOI	WER
2.0	71	128	0.24	0.41	0.99	1.30	0.53	0.99	1.01	0.51	0.97	1.66	0.74	0.86	
			-	-	-	441.7%	29.3%	0.0%	320.8	24.4%	2.0%	591.7%	80.5%	13.1%	
3.0	98	85	0.51	0.56	0.96	1.58	0.66	0.95	1.28	0.67	0.91	2.11	0.84	0.59	
			-	-	-	209.8%	17.9%	1.0%	151.0%	19.6%	5.2%	313.7%	50.0%	38.5%	
4.0	159	64	0.98	0.74	0.76	1.96	0.80	0.67	1.86	0.82	0.61	2.59	0.92	0.36	
			-	-	-	100.0%	8.1%	11.8%	89.8%	10.8%	19.7%	164.3%	24.3%	52.6%	
5.0	425	51	1.94	0.90	0.36	2.61	0.92	0.34	2.75	0.93	0.35	3.23	0.97	0.26	
			-	-	-	34.5%	2.2%	5.6%	41.8%	3.3%	2.8%	66.5%	7.8%	27.8%	

- Interpolation improves both quality and intelligibility
- ISR+MSE performs poorer than interpolation
- ISR+PL is further improved from the interpolation
 - PESQ: 66.5% \rightarrow 591.7%
 - STOI: 7.8% \rightarrow 80.5%
 - WER: 13.1% \rightarrow 52.6%

Conclusion

- An ISR system that improves intermittent sensing data
- A simple 3-step architecture
 - Acoustic feature-preserved interpolation
 - Adoption of perceptual loss
 - Combination policy addresses the missing feature issue
- Significant improvements in quality, intelligibility, and classification accuracy