


Lars Thieling, Daniel Wilhelm, Peter Jax

## 1 Introduction

**Problem:** Estimate phase  $\phi$  from given magnitude spectrum  $M$  such that a consistent time signal is achieved via inverse short-time Fourier transform (ISTFT)



### Applications:

- Speech enhancement and speech separation
- Speech synthesis and voice conversion

## 3 Phase Derivatives Estimation

- Train two equally structured DNNs using combined loss:

$$\mathcal{L}_{\text{total}} = \mathcal{L}(\Delta\hat{\psi}_{\text{if}}) + \mathcal{L}(\Delta\hat{\psi}_{\text{gd}})$$


- $\mathcal{L}$  should consider  $2\pi$  ambiguity and have a limited solution space

### Novelty I - regularized cosine loss function:

$$\mathcal{L}_{\text{reg}}(\Delta\hat{\psi}) := \sum_{k,m} -\cos(\Delta\hat{\psi}(k,m)) + \lambda \cdot (\Delta\hat{\psi}(k,m))^4$$

Here:  $\lambda = \frac{1}{4000}$

- Systematic offsets occur in the calculation of  $\psi_{\text{if}}$  and  $\psi_{\text{gd}}$



- Offset in  $\psi_{\text{if}}$  can be described by the shift theorem of the DFT:

$$x(n-S) \leftrightarrow X(k) \cdot e^{-j\frac{2\pi}{N}kS}$$


- Systematic shift in  $\psi_{\text{gd}}$  can be observed empirically

### Novelty II - shift correction:

$$\begin{aligned}\psi_{\text{if}}^*(k,m) &= \mathcal{W}\left(\psi_{\text{if}}(k,m) - \frac{\pi}{2}k\right) \\ \psi_{\text{gd}}^*(k,m) &= \mathcal{W}\left(\psi_{\text{gd}}(k,m) + \pi\right)\end{aligned}$$

DFT: discrete Fourier transform  
 $S = \frac{N}{4}$ : window shift  
 $N$ : DFT size  
 $\mathcal{W}(\cdot)$ : wrapping operator

## 2 System Overview



### Two-stage phase reconstruction system (similar to [1]):

1. Use deep neural networks (DNNs) to estimate phase derivatives 3

$$\begin{aligned}\psi_{\text{if}}(k,m) &:= \Delta_t\phi(k,m) = \phi(k,m) - \phi(k,m-1) & k: \text{freq. bin index} \\ \psi_{\text{gd}}(k,m) &:= \Delta_f\phi(k,m) = \phi(k,m) - \phi(k-1,m) & m: \text{frame index}\end{aligned}$$

2. Reconstruct phase from its estimated derivatives 4

### Proposed improvements:

- I. A novel regularized cosine loss function
- II. Shift correction (SC) as a pre-processing step
- III. A novel phase reconstruction method

## 4 Phase Reconstruction Method

- Combine  $\hat{\psi}_{\text{if}}$  and  $\hat{\psi}_{\text{gd}}$  such that a consistent  $\hat{\phi}$  is achieved


- **Novelty III** - averaging of weighted estimates  $\varphi_p$  from  $P$  paths:

$$\hat{\phi}(k,m) = \angle \sum_{p=1}^P \alpha_p(k,m) \cdot e^{j\varphi_p(k,m)}$$

with estimation quality indicators  $\alpha_p$ :

$$\begin{aligned}\alpha_1(k,m) &= M(k-1,m) \\ \alpha_2(k,m) &= M(k,m-1) \\ \alpha_3(k,m) &= \min_{l=\{-1,0\}} M(k+1,m+l)\end{aligned}$$


- Polar histograms of path error  $\varphi_p(k,m) - \phi(k,m)$  demonstrate suitability of chosen weights



## 5 Evaluation

Validation accuracies of different DNN configurations during training:

### Without SC



→  $\mathcal{L}_{\text{MSE}}$  is inappropriate


→ SC increases accuracy in first epoch

Accuracy: • Mean cosine error  
Dataset: • 18.5 hours training data  
• 3.5 hours validation data  
• 16 kHz sample rate  
DNN: • 3 hidden layers à 1024 units  
• Normalized log magnitude of frames at  $\pm 2, \pm 1$  and 0 as input features  
• Varying activation function: sigmoid, tanh, ReLU, LeakyReLU, gated linear, gated tanh  
STFT: • 640 samples Hann window  
• 160 samples window shift  
• 640 DFT size

→ SC stabilizes training of GD

→ SC stabilizes against hyperparameter variations

Results after phase reconstruction using different methods:



<http://iks.rwth-aachen.de/qr/icassp2021-rpu>

## 6 Conclusion

- Proposed novelties significantly improve phase reconstruction system
- Novelty I - regularized cosine loss function stabilizes training
- Novelty II - shift correction further stabilizes and accelerates training
- Novelty III - phase reconstruction method outperforms reference algorithms

## References

- [1] Y. Masuyama, K. Yatabe, Y. Koizumi, Y. Oikawa, and N. Harada, "Phase reconstruction based on recurrent phase unwrapping with deep neural networks," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2020.