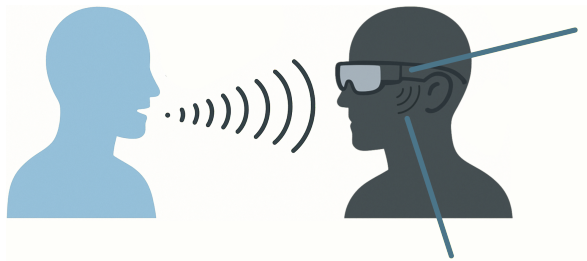# Efficient Neural and Numerical Methods for High-Quality Online Speech Spectrogram Inversion via Gradient Theorem

Andres Fernandez, Juan Azcarreta, Çağdaş Bilen, Jesus M. Alvarez

`a.fernandez@uni-tuebingen.de, jsmalvarez@meta.com`

Meta Reality Labs

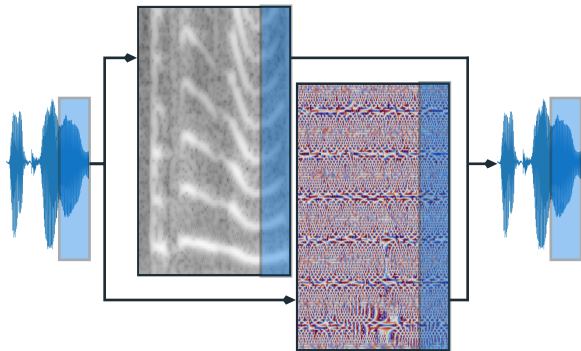INTERSPEECH 2025

17-21 AUGUST — ROTTERDAM

∞ Meta

# Motivation

Wishlist:

- ▶ Minimal **latency**

- ▶ Minimal **energy/arithmetic/memory**

- ▶ Good **quality** and **clarity**

Tasks:
{
Speech enhancement
Live translation
. . .
}

Here: **Spectrogram-based methods**

# Waveforms & Spectrograms



► Spectrograms → nice!

► Phases → messy! (irregular & $2\pi$-periodic)

► Missing phase → Inverse STFT not trivial

► Modified spectrograms may be inconsistent

Here: **Real-Time Spectrogram Inversion (RSI)**

# Efficient Neural and Numerical Methods

∞ Meta

…for **real-time** spectrogram inversion. Improvements on previous 2-stage work:

|  | GL | RTISI | SPSI | GT+DL | Ours |
|---|---|---|---|---|---|
| **Low-latency** | ✗ | ✓ | ✓ | ✓ | ✓ |
| **Low-compute** | ✗ | ✗ | ✓ | ✗ | ✓ |
| **High-quality** | ✗ | ✗ | ✗ | ✓ | ✓ |

► ∼30x smaller/faster causal CNN

► Extra 2x at cost of 1 hop in latency

► Linear-complexity least-squares solver

**Background**

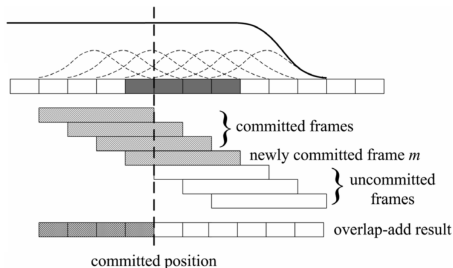# Consistency and Correctness

### Griffin-Lim (Griffin et al. 1983)



Alternating GL projections (from Peer et al. 2022)

### RTISI (Beauregard et al. 2005)



GL on the current frame alone (from Zhu et al. 2007)

▶ Initialize: $\hat{Y} \leftarrow (|Y|, \phi)$ for some phase $\phi$

▶ Consistency: $\hat{Y} \leftarrow \text{STFT} \circ \text{iSTFT} \circ \hat{Y}$

▶ Correctness: $\hat{Y} \leftarrow |Y| \frac{\hat{Y}}{|\hat{Y}|}$

▶ Recovery: $\hat{y} \leftarrow \text{iSTFT} \circ \hat{Y}$

**Real-time, but...**

▶ Requires iterations

▶ Artifacts

# Single-Pass Spectrogram Inversion

**SPSI (Beauregard et al. 2015)**

- ▶ Assume **Instantaneous Frequency**
- ▶ Initialize frame: $\hat{Y}_\tau \leftarrow (|Y|_\tau, \phi_\tau)$
- ▶ Inst. Freq.: $\omega \leftarrow$ spectral peaks in $\hat{Y}_\tau$
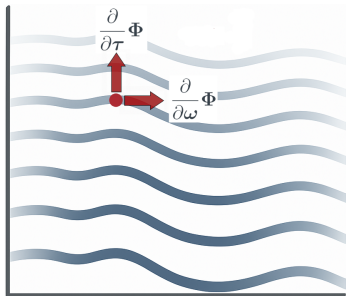- ▶ Propagate phase: $\phi_{\tau+1} \leftarrow \phi_\tau + \partial\tau \cdot \omega$

**Iteration-free**, but strong assumption $\rightarrow$ **artifacts**

# Increasing Quality with Better Assumptions

∞ Meta

**Gradient Theorem (Portnoff 1979):** for Gaussian STFT window $\varphi_\lambda(t) := e^{-\pi \frac{t^2}{\lambda}}$,

$$\frac{\partial}{\partial \omega} \operatorname{Arg}(Y_{y,\varphi_\lambda}(\omega, t)) = -\lambda \frac{\partial}{\partial t} \log |Y_{y,\varphi_\lambda}(\omega, t)|$$

$$\frac{\partial}{\partial t} \operatorname{Arg}(Y_{y,\varphi_\lambda}(\omega, t)) = \frac{1}{\lambda} \frac{\partial}{\partial \omega} \log |Y_{y,\varphi_\lambda}(\omega, t)| + 2\pi\omega$$

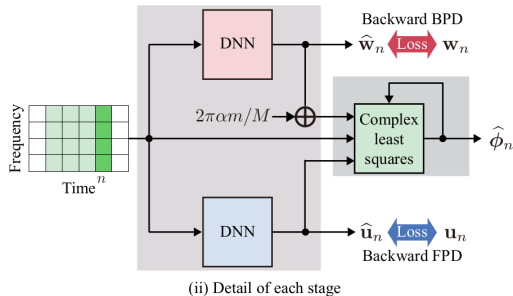Efficient numerical integration via **RTPGHI** algorithm (Průša et al. 2016)



**Powerful!**

▶ Minimal latency, $\partial$ is local

▶ No assumptions on $y$, only $\varphi$

▶ Still, error due to discretization and non-Gaussian $\varphi$

# Two-Stage Framework with Deep Learning

$\infty$ Meta

**Two stages (Masuyama et al. 2023):**

- ▶ Predict $\partial\Phi$ from $\partial|\mathbf{Y}|$ using DL
- ▶ Reconstruct $\Phi$ from $\partial\Phi$ via complex least-squares



(ii) Detail of each stage

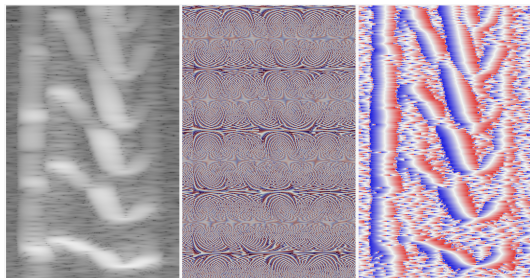Two-stage GT+DL framework from Masuyama et al. 2023

**Complex least-squares:**

$$\mathbf{z}^{(\natural)} = \arg\min_{\mathbf{z}} \underbrace{\|\mathbf{z} - (\mathbf{Y}[\omega, \tau_{\text{-}1}] \odot \mathfrak{v}_\tau)\|_{\mathbf{\Lambda}}^2}_{\tau\text{-term}} + \underbrace{\|\mathbf{D}_\tau \mathbf{z}\|_{\mathbf{\Gamma}}^2}_{\omega\text{-term}}$$

$$= \underbrace{(\mathbf{\Lambda} + \mathbf{D}_\tau^{\mathsf{H}} \mathbf{\Gamma} \mathbf{D}_\tau)^{-1}}_{A} \underbrace{\mathbf{\Lambda}(\mathbf{Y}[\omega, \tau_{\text{-}1}] \odot \mathfrak{v}_\tau)}_{b}$$

- ▶ $\mathfrak{v}$: transition from $\tau - 1$ to $\tau$

- ▶ $\mathbf{D}\mathbf{z}$: transition from $\omega$ to $\omega + 1$

- ▶ Weights $\mathbf{\Lambda}, \mathbf{\Gamma}$ to ignore small magnitudes

- ▶ Linear solver $\mathbf{z}^{(\natural)} = A^{-1}b$ for frame $\tau$

- ▶ Desired phase is $\mathrm{Arg}(\mathbf{z}^{(\natural)})$

(Recall the GT:)

$$\frac{\partial}{\partial\omega}\mathrm{Arg}(\mathsf{Y}_{y,\varphi_\lambda}(\omega,t)) = -\lambda\frac{\partial}{\partial t}\log|\mathsf{Y}_{y,\varphi_\lambda}(\omega,t)|$$

$$\frac{\partial}{\partial t}\mathrm{Arg}(\mathsf{Y}_{y,\varphi_\lambda}(\omega,t)) = \frac{1}{\lambda}\frac{\partial}{\partial\omega}\log|\mathsf{Y}_{y,\varphi_\lambda}(\omega,t)| + 2\pi\omega$$



$\log|\boldsymbol{Y}|$      $\mathrm{Arg}(\boldsymbol{Y})$      $\frac{\partial}{\partial\boldsymbol{\omega}}\mathrm{Arg}(\boldsymbol{Y})$

**Two main issues → Solutions!**

► Irregularity → **Train on derivatives!**
  ► Takamichi et al. 2018; Takamichi et al. 2020; Thieling et al. 2021; Thien et al. 2023

► $2\pi$ periodicity → **Von-Mises Loss!**
  ► $-\sum_\omega\sum_\tau\cos(\boldsymbol{X}[\omega,\tau]-\hat{\boldsymbol{X}}[\omega,\tau])$
  ► Takamichi et al. 2018; Thien et al. 2021

# Increased computation in Masuyama et al. 2023

**∞** Meta

**DNN:** 6 248k params
7.95 GMAC/s



Log-magnitude

Mean subtraction

FreqConv

FreqGatedConv

FreqConv
FreqConv
Sigmoid
⊗

FreqGatedConv
×2
⊕

FreqGatedConv
×2
⊕

FreqConv

Phase difference
(BPD or FPD)

**Complex Least-Squares:** Solving $\boldsymbol{z} = \boldsymbol{A}^{-1}\boldsymbol{b}$

$$\boldsymbol{z}_0^{(\natural)} = \underbrace{(\boldsymbol{\Lambda}_{\tau_0} + \boldsymbol{D}_{\tau_0}\boldsymbol{\Gamma}_{\tau_0}\boldsymbol{D}_{\tau_0})^{-1}}_{\boldsymbol{A}} \underbrace{\boldsymbol{\Lambda}_{\tau_0}(\boldsymbol{Y}[\omega, \tau_{\text{-}1}] \odot \mathfrak{v}_{\tau_0})}_{\boldsymbol{b}}$$

Solving $\boldsymbol{z} = \boldsymbol{A}^{-1}\boldsymbol{b}$:

▶ Memory: $\mathcal{O}(\text{L}^2)$ for STFT window of size $2\text{L}$

▶ Naive inversion of $\boldsymbol{A}$ is $\mathcal{O}(\text{L}^3)$

▶ Iterative solvers: $(\kappa(\text{L}+1)^2)$ for $\kappa$ iterations (Demmel 1997)

▶ Performed for every frame

**Very high quality**, but at **increased cost**

# Contributions

# Faster and Smaller First Stage

a) Overall Architecture

Log-Spectrogram $1\times F\times T$ → Stem → $10\times F\times T$ → Body → $10\times F\times T$ → ⊕ → $20\times F\times T$ → Head → BPD $1\times F\times T$, FPD $1\times F\times T$

b) Stem

c) Body and Head

- ▶ Cheaper, FFW layers (`BN`, `Conv1x1`, `LReLU`)

- ▶ Less residual and gated convs

- ▶ Joint FPD and BPD

- ▶ Training: Adam with CosineWR schedule

**Faster and smaller:**
- ▶ Params: 248k → 8.5k (∼30×)
- ▶ GMAC/s: 7.95 → 0.27 (∼30×)
- ▶ 2x faster, +1hop latency (★)

# Linear-Complexity Second Stage

Recall: solving $\boldsymbol{z} = \boldsymbol{A}^{-1}\boldsymbol{b}$ has complexity $\sim \mathcal{O}(\kappa \cdot L^2)$:

$$\boldsymbol{z} = \underbrace{(\boldsymbol{\Lambda} + \boldsymbol{D}^{\mathsf{H}}\boldsymbol{\Gamma}\boldsymbol{D})^{-1}}_{\boldsymbol{A}} \underbrace{\boldsymbol{\Lambda}(\boldsymbol{Y} \odot \boldsymbol{\mathfrak{v}})}_{\boldsymbol{b}}$$

Observation: $\boldsymbol{A}$ is PSD and tridiagonal!

$$\boldsymbol{D}^{\mathsf{H}}\boldsymbol{\Gamma}\boldsymbol{D} = \sum_{l=1}^{L} \gamma_l (\bar{\mathsf{d}}_l \boldsymbol{e}_l + \boldsymbol{e}_{l+1})(\mathsf{d}_l \boldsymbol{e}_l + \boldsymbol{e}_{l+1})^{\mathsf{T}}$$

$$= \sum_{l=1}^{L} \gamma_l \big( |\mathsf{d}_l|^2 \underbrace{\boldsymbol{e}_l \boldsymbol{e}_l^{\mathsf{T}}}_{\text{diag.}} + \underbrace{\boldsymbol{e}_{l+1} \boldsymbol{e}_{l+1}^{\mathsf{T}}}_{\text{diag.}} \big) + \sum_{l=1}^{L} \gamma_l \mathsf{d}_l \underbrace{\boldsymbol{e}_{l+1} \boldsymbol{e}_l^{\mathsf{T}}}_{\text{subdiag.}} + \sum_{l=1}^{L} \gamma_l \bar{\mathsf{d}}_l \underbrace{\boldsymbol{e}_l \boldsymbol{e}_{l+1}^{\mathsf{T}}}_{\text{superdiag.}}$$



Thomas' Algorithm $\rightarrow \mathcal{O}(L)$ memory and arithmetic!

# Retaining High Quality

## Intelligibility & Quality



ESTOI ↑    PESQ (WB) ↑

## More results & samples



ESTOI ↑    PESQ (WB) ↑

▶ Inversion of LibriSpeech consistent spectrograms

▶ Consistently good results on both axes

▶ Strided version also competitive

▶ Variation study supports design choices

# Thank you!

**ᐒ Meta**

**Conclusion:**

► Low latency and high quality from DL + Gradient Theorem

► Tiny causal CNN for joint BPD/FPD
  ► 2x inference at 1-hop extra latency

► Linear-complexity LSTSQ phase recovery

**Future work:**

► Subjective metrics

► Inconsistent/modified spectrograms

► Noisy phase as prior during inference

► Differentiable second stage
  ► $\Lambda, \Gamma$ as $\ell_2$ regularizers for DNN

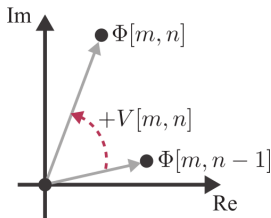**Jesus M. Alvarez**
Meta RL (Spain)

**Juan Azcarreta**
Meta RL (UK)

**Çağdaş Bilen**
Meta RL (UK)

# References

Beauregard, Gerald T., Mithila Harish, and Lonce Wyse (2015). "Single Pass Spectrogram Inversion". In: IEEE International Conference on Digital Signal Processing (DSP).

Beauregard, Gerald T, Xinglei Zhu, and Lonce Wyse (2005). "An efficient algorithm for real-time spectrogram inversion". In: DAFx.

Demmel, James W. (1997). Applied Numerical Linear Algebra. Society for Industrial and Applied Mathematics.

Griffin, D. and Jae Lim (1983). "Signal estimation from modified short-time Fourier transform". In: ICASSP.

Masuyama, Yoshiki et al. (2023). "Online Phase Reconstruction via DNN-Based Phase Differences Estimation". In: TASLP 31, pp. 163–176.

Peer, Tal, Simon Welker, and Timo Gerkmann (2022). "Beyond Griffin-Lim: Improved Iterative Phase Retrieval for Speech". In: 2022 International Workshop on Acoustic Signal Enhancement (IWAENC), pp. 1–5. DOI: 10.1109/IWAENC53105.2022.9914686.

Portnoff, M. (1979). "Magnitude-phase relationships for short-time Fourier transforms based on Gaussian analysis windows". In: ICASSP.

Průša, Zdeněk and Peter L. Søndergaard (2016). "Real-Time Spectrogram Inversion Using Phase Gradient Heap Integration". In: DAFx.

Takamichi, Shinnosuke et al. (2018). "Phase Reconstruction from Amplitude Spectrograms Based on Von-Mises-Distribution Deep Neural Network". In: IWAENC.

– (Apr. 2020). "Phase reconstruction from amplitude spectrograms based on directional-statistics deep neural networks". In: Elsevier Signal Processing 169.C.

Thieling, Lars, Daniel Wilhelm, and Peter Jax (2021). "Recurrent Phase Reconstruction Using Estimated Phase Derivatives from Deep Neural Networks". In: ICASSP.

Thien, Nguyen Binh et al. (2021). "Two-stage phase reconstruction using DNN and von Mises distribution-based maximum likelihood". In: APSIPA ASC.

– (2023). "Inter-Frequency Phase Difference for Phase Reconstruction Using Deep Neural Networks and Maximum Likelihood". In: TASLP 31.

Zhu, Xinglei, Gerald T. Beauregard, and Lonce L. Wyse (2007). "Real-Time Signal Estimation From Modified Short-Time Fourier Transform Magnitude Spectra". In: IEEE Transactions on Audio, Speech, and Language Processing 15.5, pp. 1645–1653. DOI: 10.1109/TASL.2007.899236.

$$|\mathfrak{u}_{\tau_0}| := \frac{|\boldsymbol{Y}[\omega,\tau_0]|}{|\boldsymbol{Y}[\omega\text{-}1,\tau_0]|}, \qquad \mathrm{Arg}(\mathfrak{u}_{\tau_0}) := \mathrm{Arg}\left(\frac{\boldsymbol{Y}[\omega,\tau_0]}{\boldsymbol{Y}[\omega\text{-}1,\tau_0]}\right) = \boldsymbol{u}_{\tau_0}$$

$$|\mathfrak{v}_{\tau_0}| := \frac{|\boldsymbol{Y}[\omega,\tau_0]|}{|\boldsymbol{Y}[\omega,\tau_{\text{-}1}]|}, \qquad \mathrm{Arg}(\mathfrak{v}_{\tau_0}) := \mathrm{Arg}\left(\frac{\boldsymbol{Y}[\omega,\tau_0]}{\boldsymbol{Y}[\omega,\tau_{\text{-}1}]}\right) = \boldsymbol{v}_{\tau_0}$$



Phase addition schematic from Masuyama et al. 2023

These ratios satisfy $\boldsymbol{Y}[\omega,\tau_0] = \boldsymbol{Y}[\omega\text{-}1,\tau_0] \odot \mathfrak{u}_{\tau_0}$ as well as $\boldsymbol{Y}[\omega,\tau_0] = \boldsymbol{Y}[\omega,\tau_{\text{-}1}] \odot \mathfrak{v}_{\tau_0}$ (assuming all $\boldsymbol{Y}[\omega,\tau] \neq 0$). This allows us to express $\boldsymbol{Y}[\omega,\tau_0]$ as the optimum of the following quadratic objective [21]:

$$\arg\min_{\boldsymbol{z}} \underbrace{\|\boldsymbol{z} - (\boldsymbol{Y}[\omega,\tau_{\text{-}1}] \odot \mathfrak{v}_{\tau_0})\|^2_{\boldsymbol{\Lambda}_{\tau_0}}}_{\tau\text{-term}} + \underbrace{\|\boldsymbol{D}_{\tau_0}\boldsymbol{z}\|^2_{\boldsymbol{\Gamma}_{\tau_0}}}_{\omega\text{-term}}$$

where $\boldsymbol{D}_{\tau_0} \in \mathbb{C}^{L \times (L+1)}$ is a matrix with $-\mathfrak{u}_{\tau_0}$ in the main diagonal, ones in the diagonal above, and zeros elsewhere. Here, $\|\boldsymbol{a}\|^2_{\boldsymbol{X}} := \boldsymbol{a}^{\mathsf{H}}\boldsymbol{X}\boldsymbol{a}$ is a weighted norm with *diagonal nonnegative* matrix $\boldsymbol{X}$, used in [21] to mitigate errors for small magnitudes. Equation 10 admits the following closed-form solution:

$$\boldsymbol{z}_0^{(\natural)} = (\boldsymbol{\Lambda}_{\tau_0} + \boldsymbol{D}_{\tau_0}^{\mathsf{H}}\boldsymbol{\Gamma}_{\tau_0}\boldsymbol{D}_{\tau_0})^{-1}\boldsymbol{\Lambda}_{\tau_0}(\boldsymbol{Y}[\omega,\tau_{\text{-}1}] \odot \mathfrak{v}_{\tau_0})$$