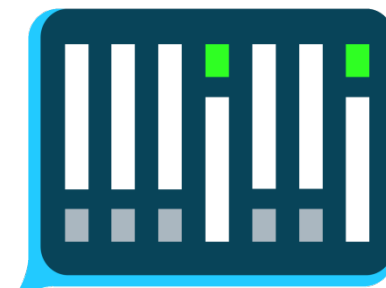# Source-filter Separation of Speech Signal in the Phase Domain

Erfan Loweimi
Jon Barker
Thomas Hain

July, 2015

# Outline

- Problems with phase spectrum

- Group delay function (GDF)

- Phase information content

- Speech signal decomposition

- Phase-based source-filter separation

- Feature extraction for ASR

- Conclusion

Erfan Loweimi

# Problems

Erfan Loweimi

# Challenges

- Historical Considerations

# Challenges

- Historical Considerations
  - Ohm's acoustic law (1843) + Helmholtz (1875)

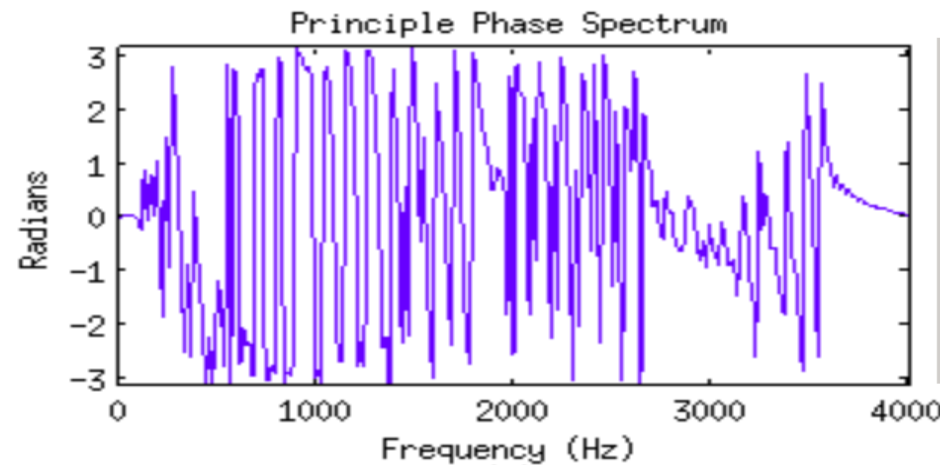Erfan Loweimi

# Challenges

- Historical Considerations
  - Ohm's acoustic law (1843) + Helmholtz (1875)
    - "the percepted quality of a tone depends solely on the *number* and *relative strength* of its partial simple tones, and not on their relative phases"

# Challenges

- Historical Considerations
  - Ohm's acoustic law (1843) + Helmholtz (1875)
    - "the percepted quality of a tone depends solely on the *number* and *relative strength* of its partial simple tones, and not on their relative phases"

  - Although some studies show that the auditory system is not totally "*phase deaf*", this law forms the status qua
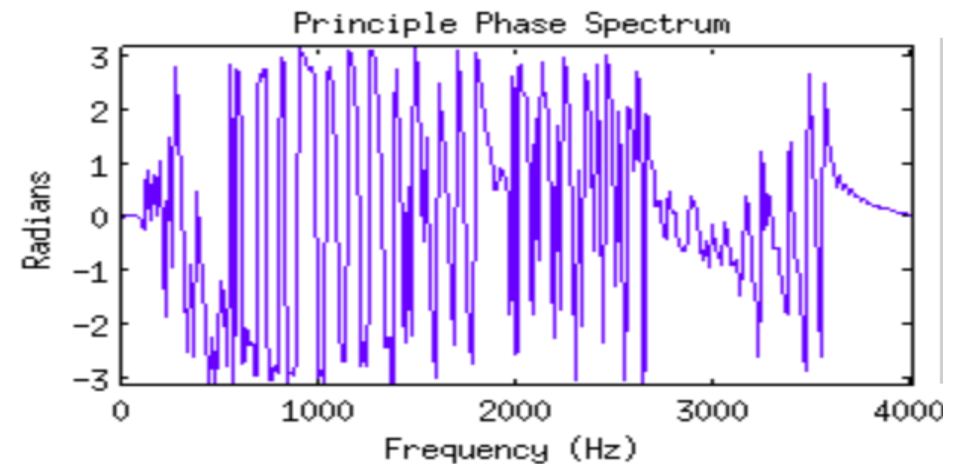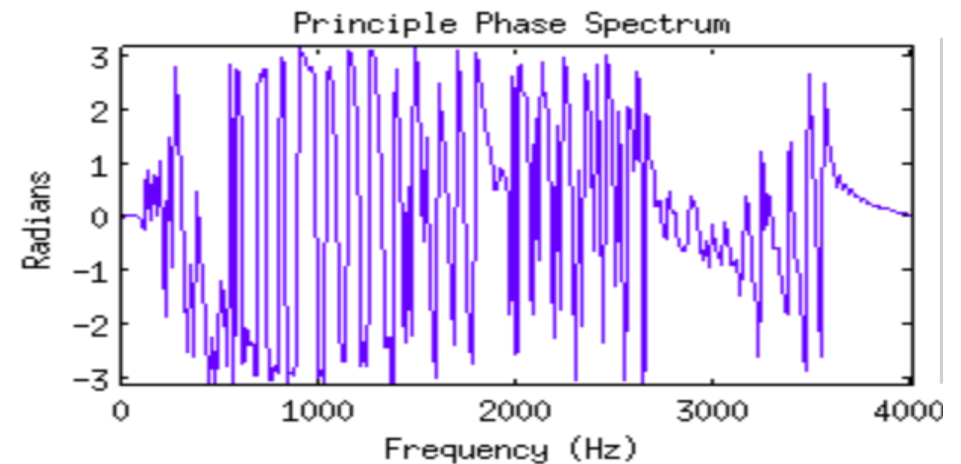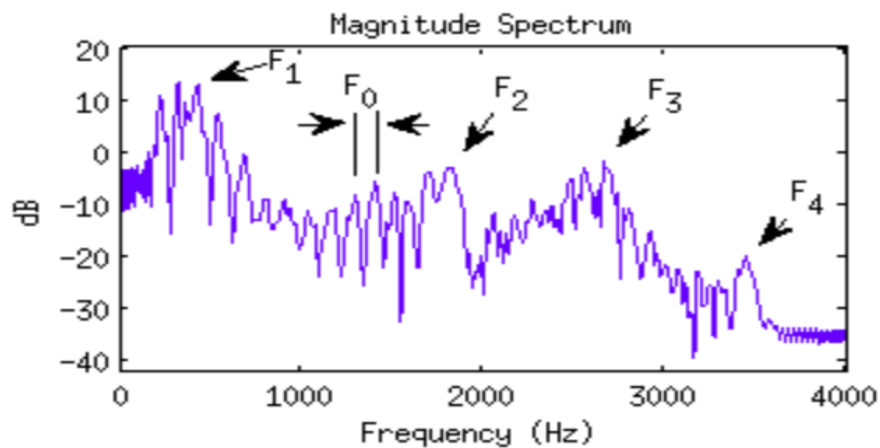
# Challenges ...

- Phase wrapping

# Challenges ...

- Phase wrapping
  - Chaotic/noise-like behaviour
  - Lacks any meaningful trend or extrema points
    - Physical interpretation
    - Mathematical modelling



Erfan Loweimi

# Challenges ...

- Phase wrapping
  - Chaotic/noise-like behaviour
  - Lacks any meaningful trend or extrema points
    - Physical interpretation
    - Mathematical modelling

# Challenges ...

- Only informative in long-term ( > x00 ms)

    – Violates stationarity assumption !

    – In short frames (~ 30 ms), it is generally believed that the phase spectrum does not contribute much to speech quality/intelligibility

Erfan Loweimi

# Group Delay Function (GDF)

- Definition

$$\tau_X(\omega) = -\frac{d}{d\omega} arg[X(\omega)] = -Im\{\frac{d}{d\omega}\log(X(\omega))\}$$

Erfan Loweimi

# Group Delay Function (GDF)

- Definition

$$\tau_X(\omega) = -\frac{d}{d\omega} arg[X(\omega)] = -Im\{\frac{d}{d\omega}\log(X(\omega))\}$$

$$\tau_X(\omega) = \frac{X_{Re}(\omega)Y_{Re}(\omega) + X_{Im}(\omega)Y_{Im}(\omega)}{|X(\omega)|^2}$$

# Group Delay Function (GDF)

- Definition

$$\tau_X(\omega) = -\frac{d}{d\omega} arg[X(\omega)] = -Im\{\frac{d}{d\omega} \log(X(\omega))\}$$
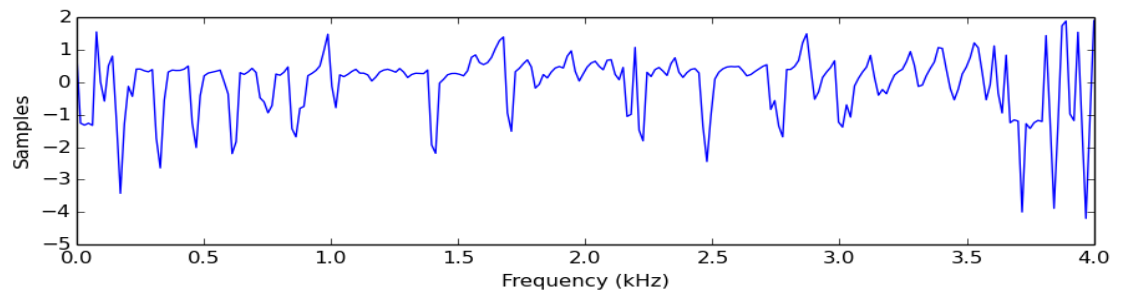
$$\tau_X(\omega) = \frac{X_{Re}(\omega)Y_{Re}(\omega) + X_{Im}(\omega)Y_{Im}(\omega)}{|X(\omega)|^2}$$

- Pros
  - Resembles the magnitude spectrum
  - High frequency resolution
  - Additive

# Group Delay Function (GDF)

- Definition

$$\tau_X(\omega) = -\frac{d}{d\omega} arg[X(\omega)] = -Im\{\frac{d}{d\omega} \log(X(\omega))\}$$

$$\tau_X(\omega) = \frac{X_{Re}(\omega)Y_{Re}(\omega) + X_{Im}(\omega)Y_{Im}(\omega)}{|X(\omega)|^2}$$

- Pros
  - Resembles the magnitude spectrum
  - High frequency resolution
  - Additive
- Cons
  - Too spiky

# Phase Information Content

- What is the information?

# Phase Information Content

- ## What is the information?

  - ### Context dependent

    - Information theory: average of uncertainty
    - Speech: lingual content, speaker ID, ...

# Phase Information Content

- What is the information?

  - Context dependent

    - Information theory: average of uncertainty

    - Speech: lingual content, speaker ID, ...

- Is phase informative?

  - ✔ From perceptual viewpoint

  - ✗ From signal processing viewpoint

# Signal Decomposition

- For any signal

$$X(\omega) = |X(\omega)| \, . \, e^{j\phi_X(\omega)}$$
$$X(\omega) = X_{MinPh}(\omega) \, . \, X_{AllPass}(\omega)$$
$$= |X_{MinPh}(\omega)|e^{j\phi_{MinPh}(\omega)} \, . \, 1e^{j\phi_{AllPass}(\omega)}$$

# Signal Decomposition

- For any signal

$$X(\omega) = |X(\omega)| \, . \, e^{j\phi_X(\omega)}$$
$$X(\omega) = X_{MinPh}(\omega) \, . \, X_{AllPass}(\omega)$$
$$= |X_{MinPh}(\omega)|e^{(j\phi_{MinPh}(\omega)+\phi_{AllPass}(\omega))}$$

$$\begin{cases} |X(\omega)| = |X_{MinPh}(\omega)| \\ \phi_X(\omega) = \phi_{MinPh}(\omega) + \phi_{AllPass}(\omega) \end{cases}$$

# Signal Decomposition

- For any signal

$$\begin{cases} |X(\omega)| = |X_{MinPh}(\omega)| \\ \phi_X(\omega) = \phi_{MinPh}(\omega) + \phi_{AllPass}(\omega) \end{cases}$$

\* Is there any relation between phase and magnitude spectra?

# Signal Decomposition

- For any signal

$$\begin{cases} |X(\omega)| = |X_{MinPh}(\omega)| \\ \phi_X(\omega) = \phi_{MinPh}(\omega) + \phi_{AllPass}(\omega) \end{cases}$$

* Is there any relation between phase and magnitude spectra?

$$|X_{MinPh}(\omega)| \quad \underleftrightarrow{\text{HiL.Tran}} \quad \phi_{MinPh}(\omega)$$

# Signal Decomposition



- For any signal

$$\begin{cases} |X(\omega)| = |X_{MinPh}(\omega)| \\ \phi_X(\omega) = \phi_{MinPh}(\omega) + \phi_{AllPass}(\omega) \end{cases}$$

\* Is there any relation between phase and magnitude spectra?

$$|X_{MinPh}(\omega)| \quad \underleftrightarrow{\text{HiL.Tran}} \quad \phi_{MinPh}(\omega)$$

# Signal Decomposition ...

- For speech ...

$$\begin{cases} |X(\omega)| = |X_{VT}(\omega)|.|X_{Exc}(\omega)| \\ \\ \end{cases}$$

# Signal Decomposition ...

- For speech ...

$$\left\{ \begin{array}{l} |X(\omega)| = |X_{VT}(\omega)|.|X_{Exc}(\omega)| = |X_{MinPh}(\omega)| \\ \\ \\ \end{array} \right.$$

# Signal Decomposition ...

- For speech

$$
\begin{cases}
|X(\omega)| = |X_{VT}(\omega)|.|X_{Exc}(\omega)| = |X_{MinPh}(\omega)| \\
|X_{MinPh}(\omega)| \quad \underrightarrow{\text{HiL.Tran}} \quad \arg[X_{MinPh}(\omega)]
\end{cases}
$$

# Signal Decomposition ...

- For speech

$$
\begin{cases}
|X(\omega)| = |X_{VT}(\omega)|.|X_{Exc}(\omega)| = |X_{MinPh}(\omega)| \\[2ex]
|X_{MinPh}(\omega)| \quad \underleftrightarrow{\text{HiL.Tran}} \quad \arg[X_{MinPh}(\omega)] \\[2ex]
arg[X_{MinPh}(\omega)] = arg[X_{VT}(\omega)] + arg[X_{Exc}(\omega)].
\end{cases}
$$

# Signal Decomposition ...

- For speech

$$\begin{cases} |X(\omega)| = |X_{VT}(\omega)|.|X_{Exc}(\omega)| = |X_{MinPh}(\omega)| \\ |X_{MinPh}(\omega)| \quad \underset{\xleftarrow{\hspace{2em}}}{\text{HiL.Tran}} \quad arg[X_{MinPh}(\omega)] \\ arg[X_{MinPh}(\omega)] = arg[X_{VT}(\omega)] + arg[X_{Exc}(\omega)]. \end{cases}$$

- Goal ...



Phase-based processing

$x[n]$

$x_{VT}[n] * x_{Exc}[n]$

$\hat{x}_{VT}[n]$

$\hat{x}_{Exc}[n]$

# Signal Decomposition ...

- For speech

$$
\begin{cases}
|X(\omega)| = |X_{VT}(\omega)|.|X_{Exc}(\omega)| = |X_{MinPh}(\omega)| \\[2mm]
|X_{MinPh}(\omega)| \quad \underset{\xleftrightarrow{\text{HiL.Tran}}}{} \quad arg[X_{MinPh}(\omega)] \\[2mm]
arg[X_{MinPh}(\omega)] = arg[X_{VT}(\omega)] + arg[X_{Exc}(\omega)].
\end{cases}
$$

- Goal ...

Phase-based processing

$$x[n]$$

$$x_{VT}[n] * x_{Exc}[n] \rightarrow arg[X_{MinPh}(\omega)] \nearrow arg[X_{VT}(\omega)] \rightarrow \hat{x}_{VT}[n]$$

$$\searrow arg[X_{Exc}(\omega)] \rightarrow \hat{x}_{Exc}[n]$$

# MinPh Component Computation

- In Frequency domain

- In Quefrency domain

# MinPh Component Computation

- In Frequency domain

$$\arg[\mathrm{X}_{MinPh}(\omega)] = Hil\{log|X_{MinPh}(\omega)|\}$$
$$= -\frac{1}{2\pi}log|X_{MinPh}(\omega)| * cot(\frac{\omega}{2})$$

- In Quefrency domain

    – Apply a proper lifter on the complex cepstrum

# MinPh Component Computation

- In Frequency domain

$$\arg[\mathrm{X}_{MinPh}(\omega)] = Hil\{log|X_{MinPh}(\omega)|\}$$
$$= -\frac{1}{2\pi}log|X_{MinPh}(\omega)| * cot(\frac{\omega}{2})$$

- In Quefrency domain

  – Apply a proper lifter on the complex cepstrum

MinPh sequence => causal cepstrum

# MinPh Component Computation

- In Frequency domain

$$\arg[\mathrm{X}_{MinPh}(\omega)] = Hil\{log|X_{MinPh}(\omega)|\}$$
$$= -\frac{1}{2\pi}log|X_{MinPh}(\omega)| * cot(\frac{\omega}{2})$$

- In Quefrency domain

  - Apply a proper lifter on the complex cepstrum



Speech is mixed-phase

# MinPhase component

$$log|X(\omega)|$$



**Hilbert Transform**

$$ARG[X(\omega)]$$



$$arg[X_{MinPh}(\omega)]$$

# MinPhase component

$$log|X(\omega)|$$



Hilbert Transform



$$ARG[X(\omega)]$$



$$arg[X_{MinPh}(\omega)]$$

# Trend/Fluctuation Analysis

# Trend/Fluctuation Analysis



$$arg[X_{MinPh}] = \ {\color{red}Trend} \ + \ {\color{green}Fluctuation}$$

# Trend/Fluctuation Analysis



$$arg[X_{MinPh}] = \textcolor{red}{Trend} + \textcolor{green}{Fluctuation}$$

# Trend/Fluctuation Analysis



$$arg[X_{MinPh}] = \textcolor{red}{Trend} + \textcolor{green}{Fluctuation}$$

Vocal Tract (Filter)

Excitation (Source)

Erfan Loweimi

# Trend/Fluctuation Separation



Fourier Transform

X          Y

Underlying assumption for having successful separation

# Trend/Fluctuation Separation



Underlying assumption for having successful separation

# Trend/Fluctuation Separation



Fourier Transform

X          Y

Underlying assumption for having successful separation

y + x

Fourier Transform

X          Y

y + x

Fourier Transform

V          E

e + v

Speech signal

Erfan Loweimi

13/18

# Phase-based Source-Filter Decomposition

# Phase-based Source-Filter Decomposition



$$\begin{cases} \hat{\tau}_{VT}(\omega) = signum(\tau_{VT}(\omega)).|\tau_{VT}(\omega)|^{\alpha} \\ signum(\tau_{VT}(\omega)) = \frac{\tau_{VT}(\omega)}{|\tau_{VT}(\omega)|} \end{cases}$$

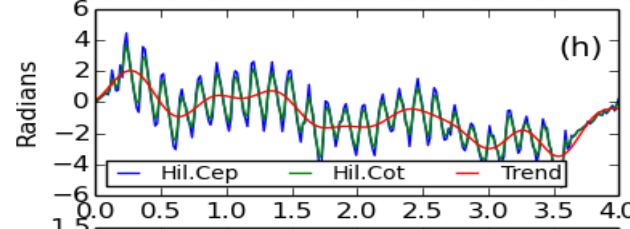# Phase-based Source-Filter Decomposition
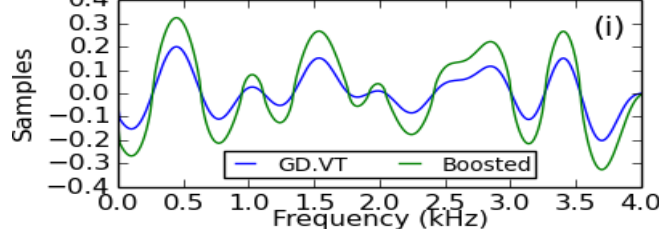


waveform

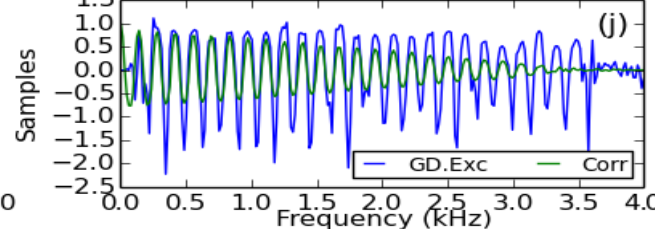Phase Spec.

Mag. Spec.

MODGDF(03)

CGDF

Product Spec.

AR+GDF

MinPh Phase

GDF.VT
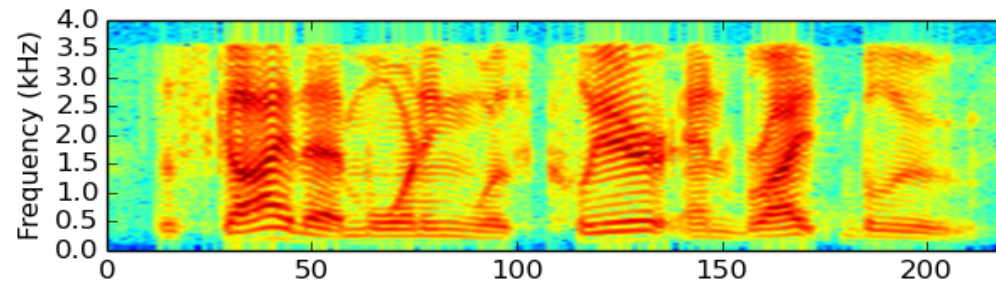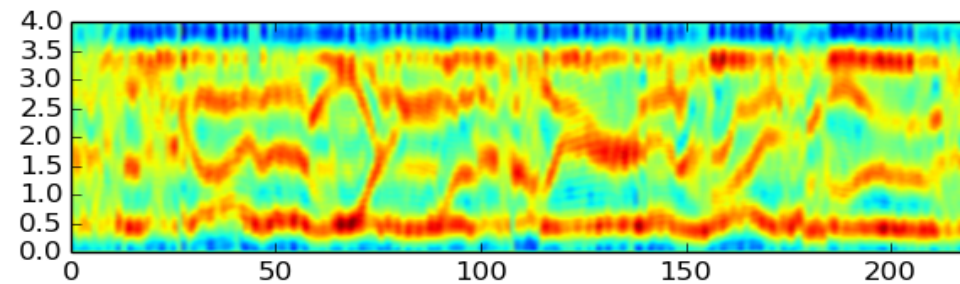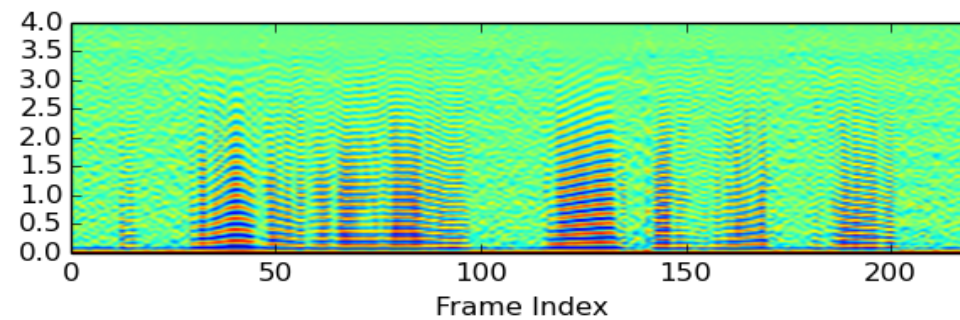
GDF.Exc

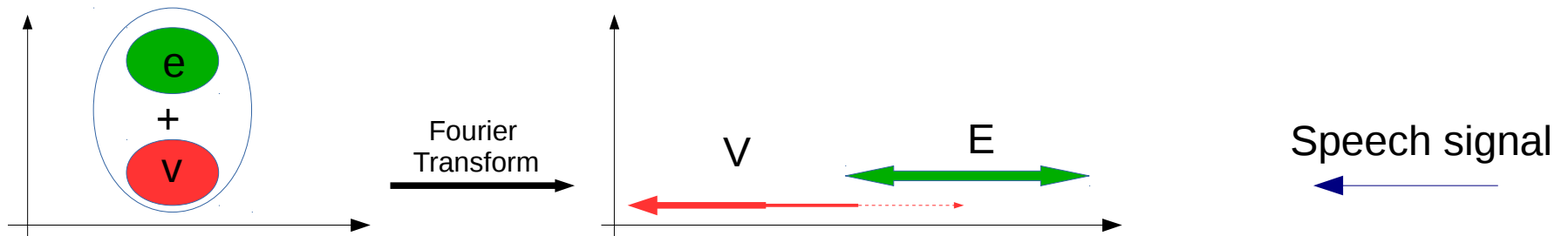# Phase-based Source-Filter Decomposition
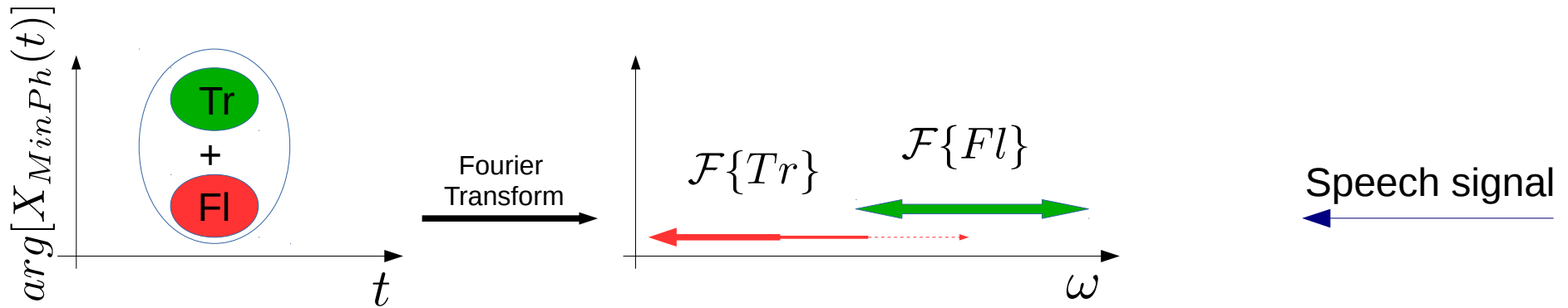
$log|X(n,\omega)|$

$\tau_{VT}(n,\omega)$

$\tau_{Exc}(n,\omega)$
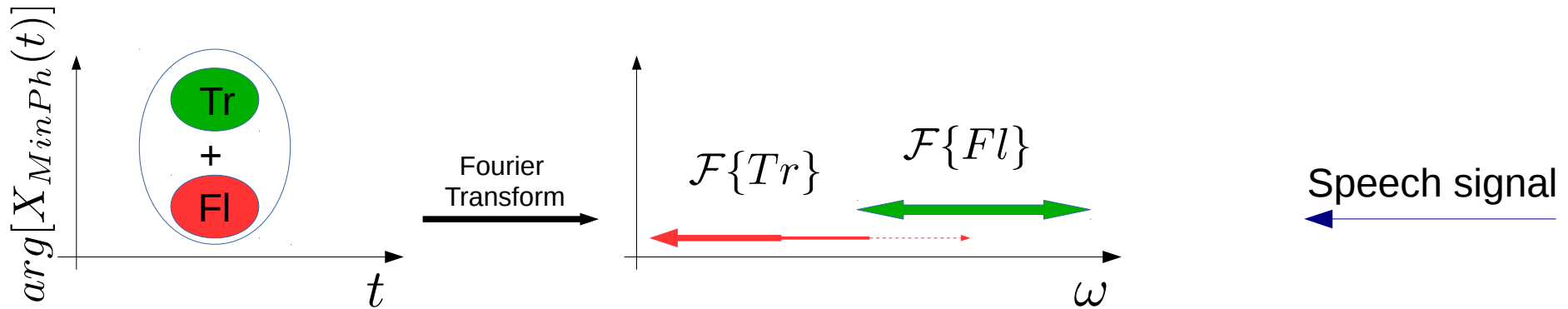
# Decomposition in GDF domain

# Decomposition in GDF domain

# Decomposition in GDF domain



$$
\begin{cases}
\tau_X(t) = -\frac{d}{dt}arg[X_{MinPh}(t)] = -\frac{d}{dt}Trend - \frac{d}{dt}Fluctuation \\
\mathcal{F}\{\tau_X(t)\} = -j\omega\mathcal{F}\{Trend\} - j\omega\mathcal{F}\{Fluctuation\}
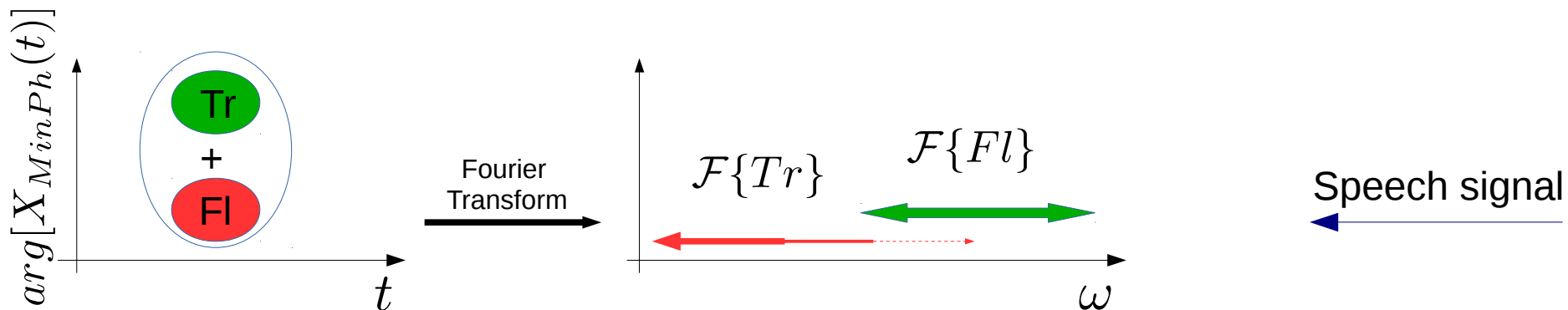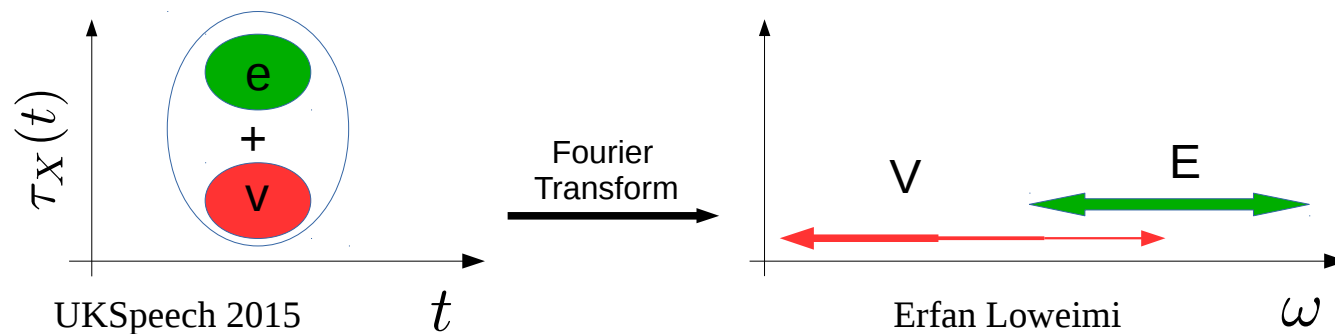\end{cases}
$$

Erfan Loweimi

# Decomposition in GDF domain



$$\begin{cases} \tau_X(t) = -\frac{d}{dt}arg[X_{MinPh}(t)] = -\frac{d}{dt}Trend - \frac{d}{dt}Fluctuation \\ \mathcal{F}\{\tau_X(t)\} = -j\omega\mathcal{F}\{Trend\} - j\omega\mathcal{F}\{Fluctuation\} \end{cases}$$

Erfan Loweimi

# Decomposition in GDF domain



$$\begin{cases} \tau_X(t) = -\frac{d}{dt}arg[X_{MinPh}(t)] = -\frac{d}{dt}Trend - \frac{d}{dt}Fluctuation \\ \mathcal{F}\{\tau_X(t)\} = -j\omega\mathcal{F}\{Trend\} - j\omega\mathcal{F}\{Fluctuation\} \end{cases}$$
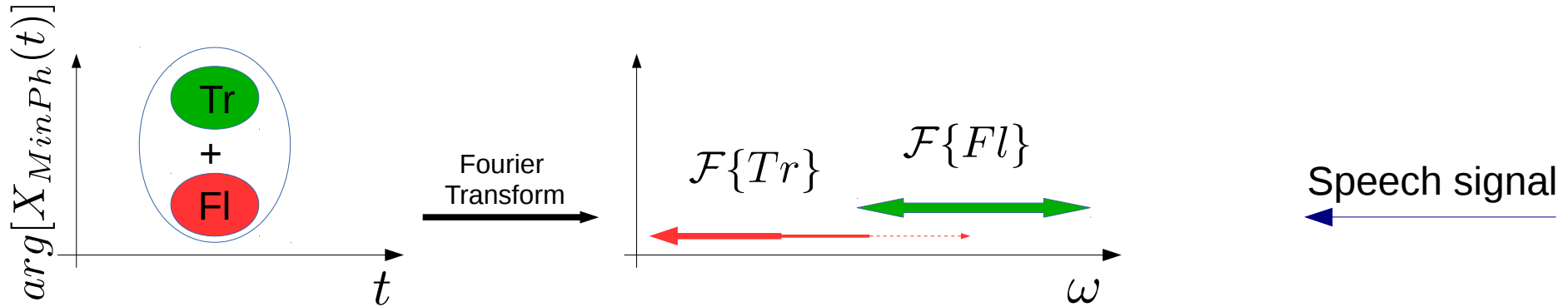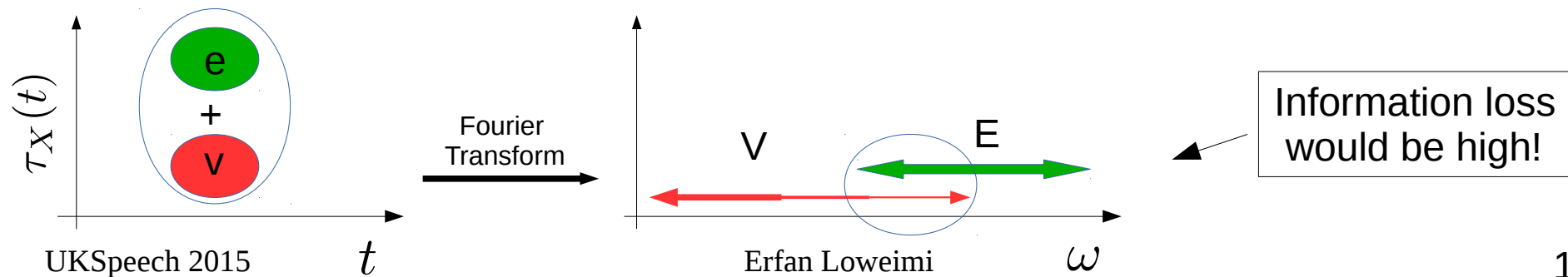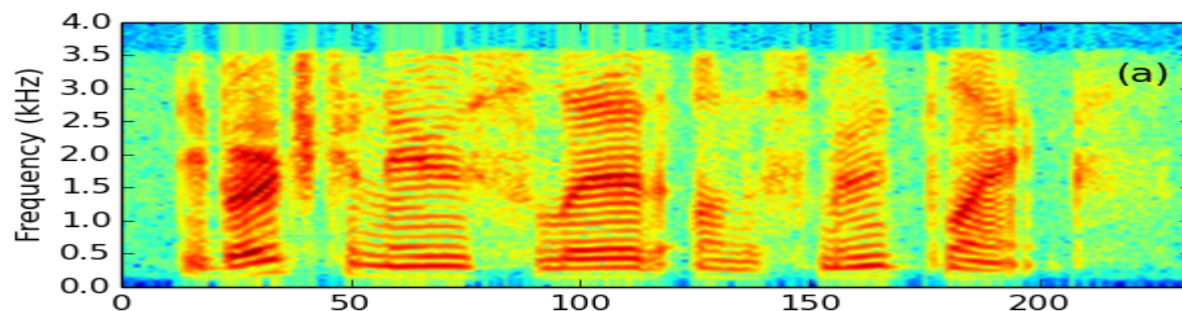


Information loss would be high!
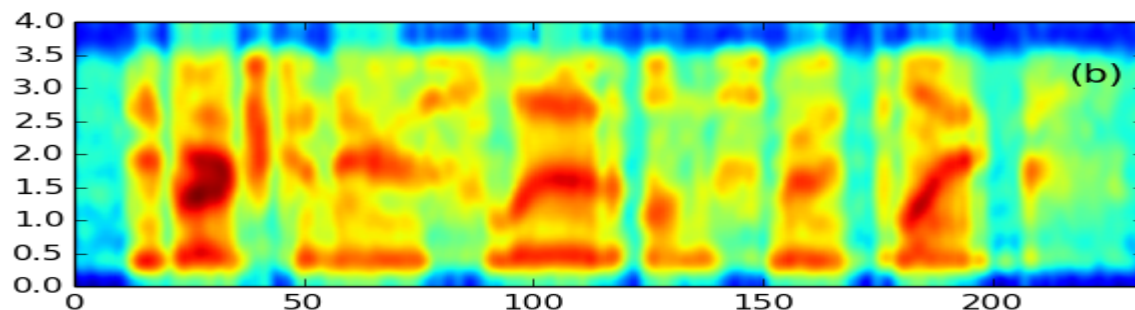
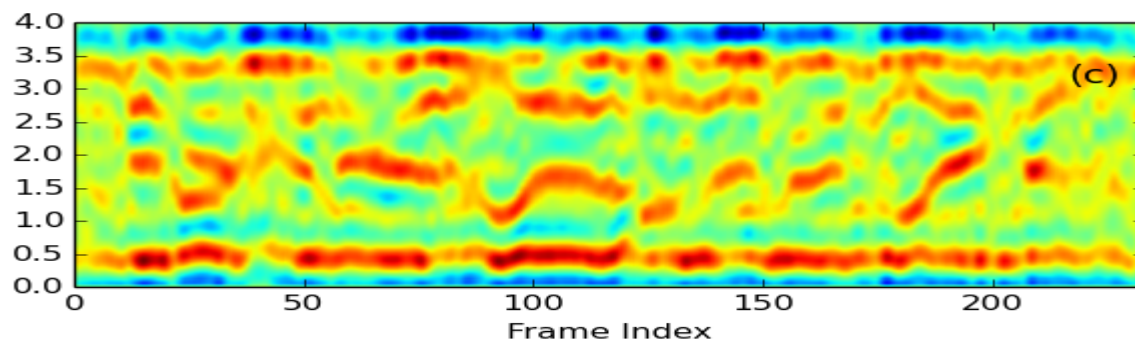# Decomposition in log-magnitude domain

$log|X(n,\omega)|$

Smoothed Spec.

$\tau_X(n,\omega)$

Erfan Loweimi

# Decomposition in log-magnitude domain

$log|X(n,\omega)|$

Smoothed Spec.

$\tau_X(n,\omega)$



Erfan Loweimi

# Feature Extraction for ASR

$$\text{i)} \quad arg[X_{VT}] \rightarrow DCT \Rightarrow PHVT$$

$$\text{ii)} \quad \tau_{VT} \rightarrow DCT \Rightarrow GDVT$$

$$\text{iii)} \quad \tau_{VT} \rightarrow MelFilterbank \rightarrow DCT \Rightarrow MFGDVT$$

$$\text{iv)} \quad \tau_{VT} \rightarrow Mel\ Filterbank \rightarrow Boost \rightarrow DCT \Rightarrow BMFGDVT$$

# Feature Extraction for ASR

| Feature | TestSet A | TestSet B | TestSet C |
|---|---|---|---|
| MFCC | 66.2 | 71.4 | 64.9 |
| PLP | 67.3 | 70.6 | 66.2 |
| PNCC | 71.2 | 72.8 | 71.5 |
| MODGDF | 64.3 | 66.4 | 59.5 |
| CGDF | 67.0 | 73.0 | 59.4 |
| PS | 66.0 | 71.2 | 64.6 |
| i) PHVT | 69.0 | 74.8 | 67.1 |
| ii) GDVT | 70.5 | 75.9 | 69.1 |
| iii) MFGDVT | 72.8 | 77.3 | 72.8 |
| iv) BMFGDVT | **73.2** | **77.4** | **73.4** |

i) $arg[X_{VT}] \to DCT \Rightarrow PHVT$

ii) $\tau_{VT} \to DCT \Rightarrow GDVT$

* Aurora 2
* Average of 0-20 dB

UKSpeec iii) $\tau_{VT} \to MelFilterbank \to DCT \Rightarrow MFGDVT$

iv) $\tau_{VT} \to Mel\ Filterbank \to Boost \to DCT \Rightarrow BMFGDVT$

# Conclusion

- This talk was about phase-based source-filter deconvolution

- Separation was done using Trend/Fluctuation analysis of the phase spectrum of the minimum-phase component of speech

- Proposed method succeeds in decomposing the speech into vocal tract and excitation components

- Extracted feature from the vocal tract component of the phase shows good robustness on Aurora 2 task

Erfan Loweimi

# That is it!

- Thanks for your attention
- Question

Erfan Loweimi