# Audio **&** Speech **Tech**nology

[4] Audio Features

---

**Lower false negative**

Robustness

**Lower false positive**

Reliability

## Audio **Features**

**Lower computational complexity**

Efficiency

---

# Audio Spectral **Envelope**

$$ASE[b] = \sum_{k \in Band[b]} P[k]$$



$$P[k] = |S[k]|^2$$

$S[k] = Spectrum$

Audio Spectrum Envelope

Power Spectrum

Band Index $b$

$ASE(b)$

$P(k)$

$f(k)$ (log$_2$-Hz)

31.25  62.5  125  250  500  **1k**  2k  4k  8k  16k  $F_s/2$

out-of-band    within-band    out-of-band

---

# Audio Spectral **Flatness**

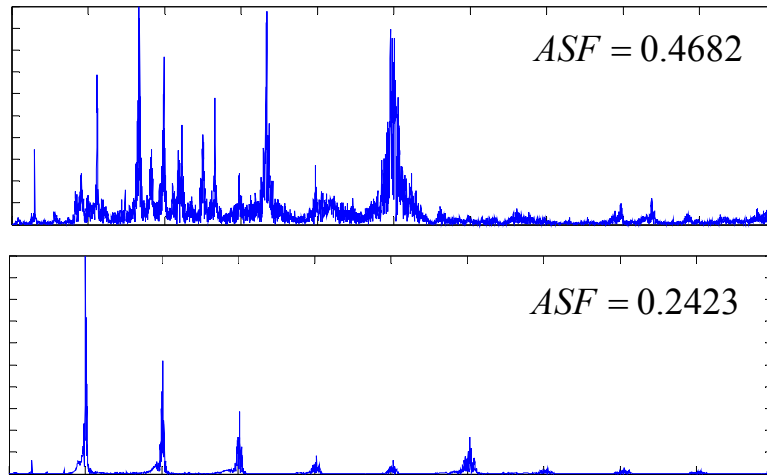$$ASF[b] = \frac{\sqrt[N_b]{\prod_{k \in Band[b]} P[k]}}{\frac{1}{N_b} \sum_{k \in Band[b]} P[k]}$$

→ **Geometric Mean**

→ **Arithmetic Mean**

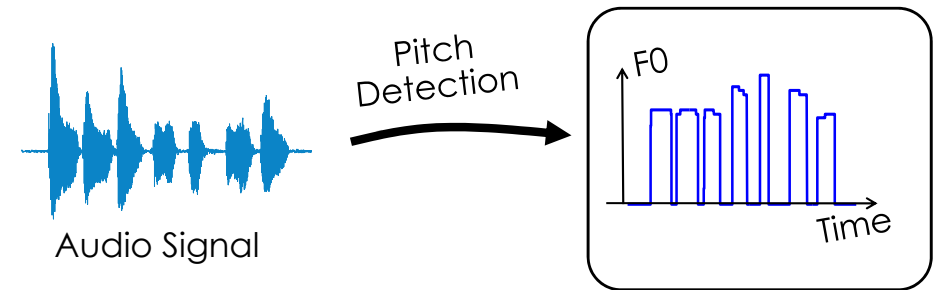$N_b = \text{No. of frequencies in each band}$

$$\sqrt[N_b]{\prod_{k \in Band[b]} P[k]} = \exp\left(\frac{1}{N_b} \sum_{k \in Band[b]} \ln P[k]\right)$$

# Audio Spectral **Flatness**

$ASF = 0.4682$

$ASF = 0.2423$

# **Fundamental** Audio **Frequency**

Pitch Detection

F0

Time

Audio Signal

# **Zero** Crossing **Rate**

$$ZCR = \frac{F_S}{2N} \sum_{n=1}^{N-1} \left| sign(s[n]) - sign(s[n-1]) \right|$$

$N = \text{No. samples in } s[i]$

$F_S = \text{Sampling Frequency}$

Waveform

Time Index

# Spectral **Centroid**

**Power Spectrum (log)**

$$SC = \frac{\sum_k f[k]P[k]}{\sum_k P[k]}$$

Centroid

-2
-3
-4
-5
-6
-7
-8

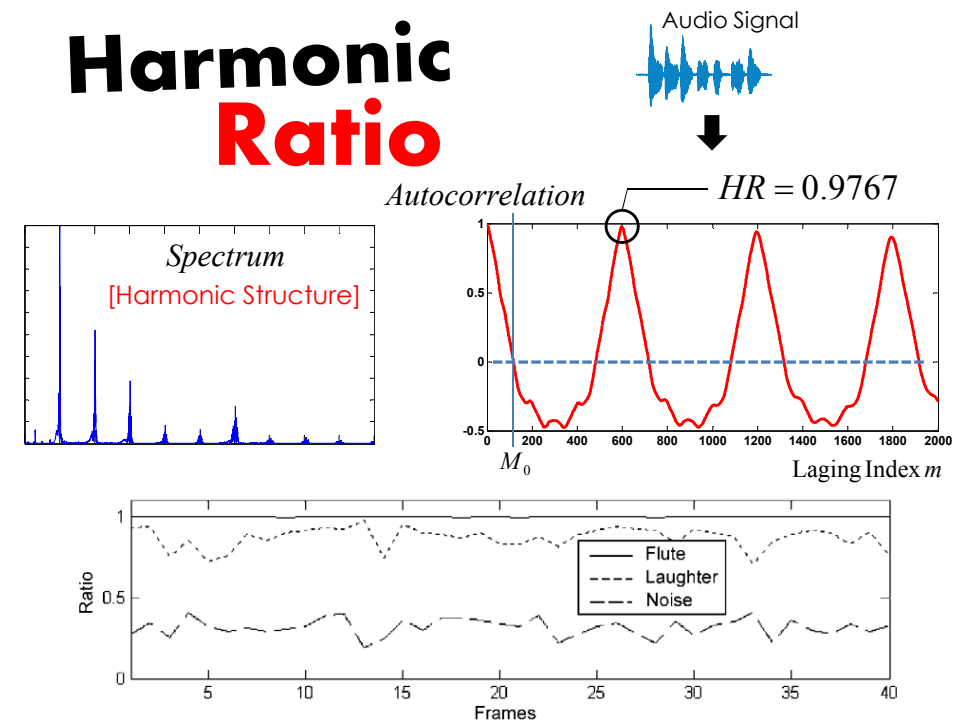50    100    150    200    250    300    350    400    450

**Frequency**

# Harmonic Ratio

$$R[m] = \frac{\sum_{n \in Frame} s[n]s[n-m]}{\sqrt{\sum_{n \in Frame} s^2[n] \sum_{n \in Frame} s^2[n-m]}} \Rightarrow HR = \max_{m \geq M_0} R[m]$$
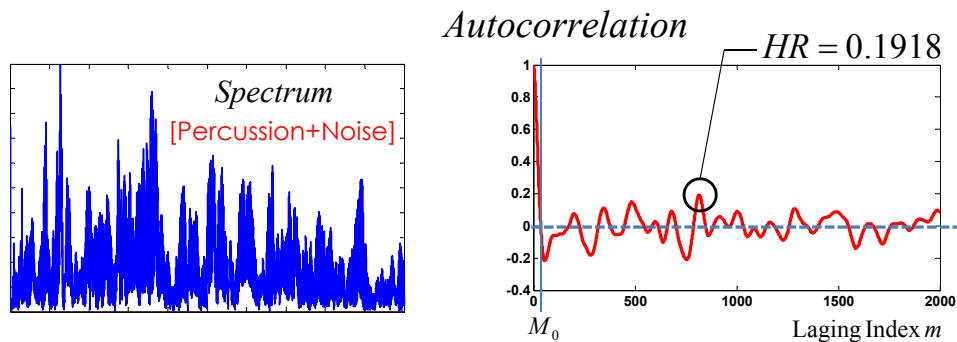
$R$ = Autocorrelation Function

$m$ = Lagging index

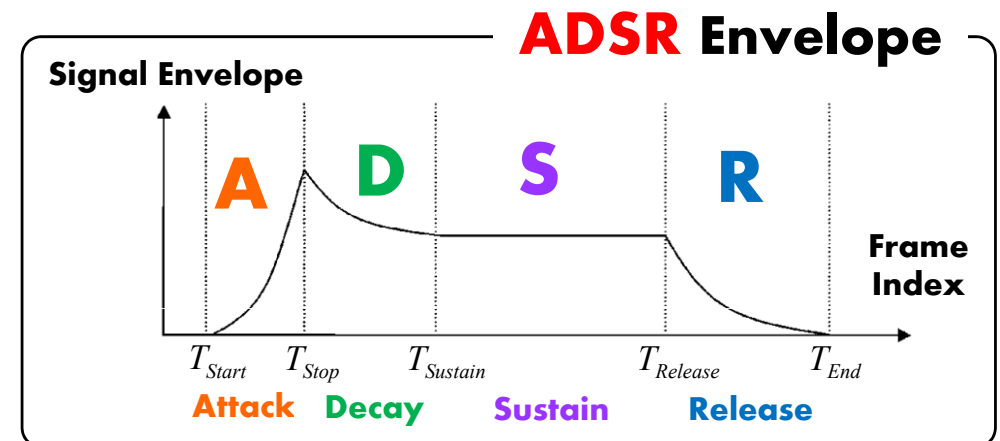$M_0$ = Position of the first zero crossing of autocorrelation $R$

# Harmonic Ratio

Audio Signal

$HR = 0.9767$

*Autocorrelation*

*Spectrum*
[Harmonic Structure]

$M_0$   Laging Index $m$

Ratio — Flute / Laughter / Noise

Frames

# Harmonic Ratio

*Autocorrelation*

$HR = 0.1918$

*Spectrum*
[Percussion+Noise]

$M_0$   Laging Index $m$

# Log Attack Time

**ADSR Envelope**

**Signal Envelope**

**A**  **D**  **S**  **R**

**Frame Index**

$T_{Start}$  $T_{Stop}$  $T_{Sustain}$  $T_{Release}$  $T_{End}$

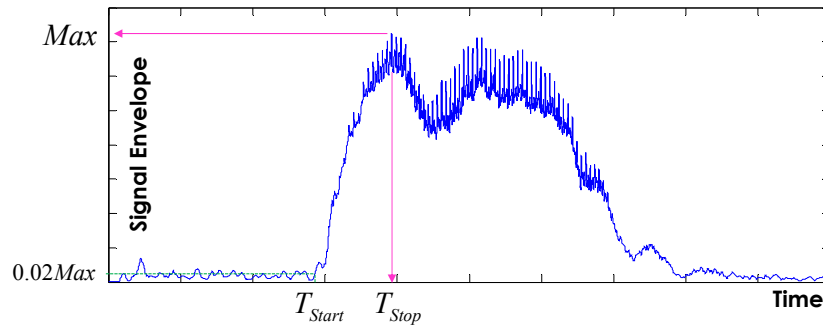**Attack**  **Decay**  **Sustain**  **Release**
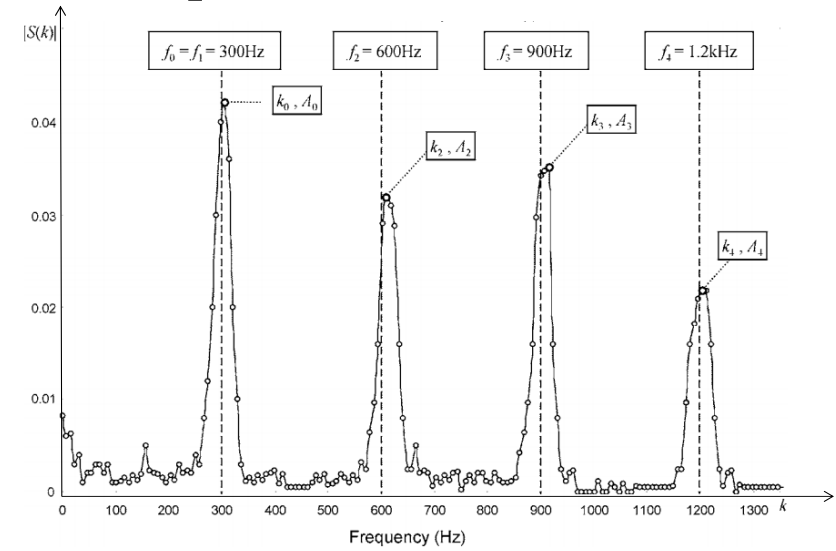
# Log Attack Time

$$LAT = \log_{10}\left(T_{Stop} - T_{Start}\right)$$

$T_{Start}$ = The time the signal envelope exceeds 2% of its maximal value

$T_{Stop}$ = The time the signal evelope reaches its maximal value

# Harmonic Spectral Centroid



$f_0 = f_1 = 300Hz$, $f_2 = 600Hz$, $f_3 = 900Hz$, $f_4 = 1.2kHz$

MPEG-7 Audio and. Beyond. Audio Content Indexing and. Retrieval. Hyoung-Gook Kim, et. al.

# Harmonic Spectral Centroid

$LHSC$ = Local Harmonic Spectral Centroid of Audio Frame

$$LHSC = \frac{\sum_{h=1}^{N_H} f_h A_h}{\sum_{h=1}^{N_H} A_h}$$

$$HSC = \frac{1}{L}\sum_{i=1}^{L} LHSC[i]$$

$f_h$ = Frequency of the $h^{th}$ harmonic

$A_h$ = Amplitude of the $h^{th}$ harmonic

$N_H$ = No. of harmonic peaks

$L$ = No. of time frames

# Harmonic Spectral Deviation

Spectral Envelope

$$SE_h = \begin{cases} (A_h + A_{h+1})/2 & ; h = 1 \\ (A_{h-1} + A_h + A_{h+1})/2 & ; h \in [2, N_H - 1] \\ (A_{h-1} + A_h)/2 & ; h = N_H \end{cases}$$

$$LHSD = \frac{\sum_{h=1}^{N_H} \left|\log_{10} A_h - \log_{10} SE_h\right|}{\sum_{h=1}^{N_H} \log_{10} A_h}$$

$$HSD = \frac{1}{L}\sum_{i=1}^{L} LHSD[i]$$

# Harmonic Spectral Spread

$$LHSS = \sqrt{\dfrac{\sum\limits_{h=1}^{N_H}(f_h - LHSC)^2 A_h^2}{\sum\limits_{h=1}^{N_H} A_h^2}}$$

$$HSS = \frac{1}{L}\sum_{i=1}^{L} LHSS[i]$$

# Harmonic Spectral Variation

$$LHSV[i] = 1 - \dfrac{\sum\limits_{h=1}^{N_H} A_h[i-1]A_h[i]}{\sqrt{\sum\limits_{h=1}^{N_H} A_h^2[i-1]}\sqrt{\sum\limits_{h=1}^{N_H} A_h^2[i]}}$$
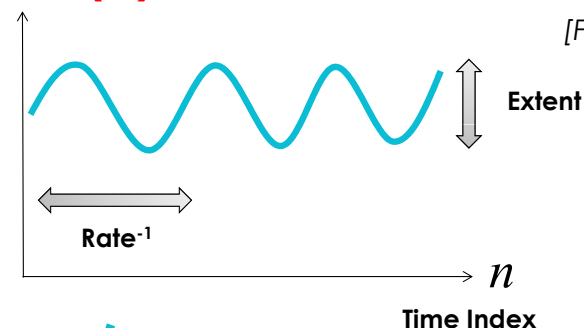
$$HSS = \frac{1}{L}\sum_{i=1}^{L} LHSS[i]$$

# Vibrato

*[Frequency Modulation]*

- Variation of pitch.
- **Extent of vibrato**
    = Amount of pitch variation
- **Rate of vibrato**
    = Speed which the pitch is varied
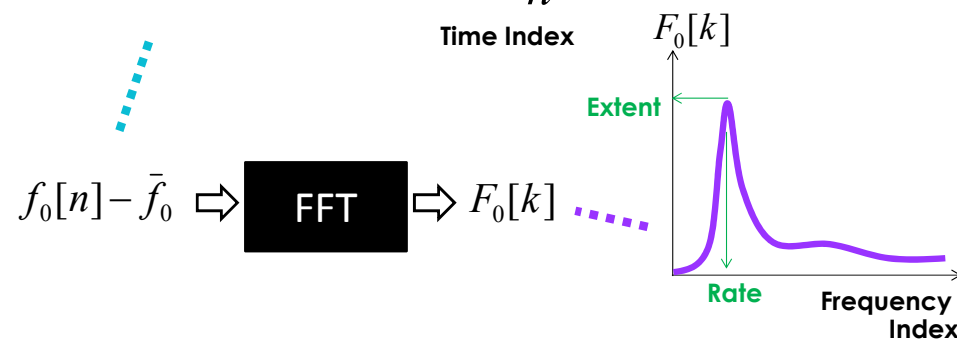- String instruments produce the FM dominant sounds

# Vibrato

*[Frequency Modulation]*
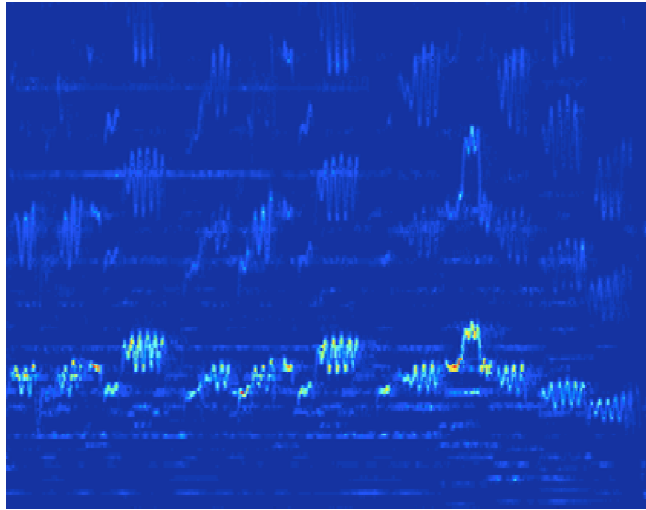
$f_0[n]$
**Pitch (f0)**

Extent

Rate⁻¹

$n$

**Time Index**

$F_0[k]$

Extent

Rate

**Frequency Index**

$f_0[n] - \bar{f}_0 \Rightarrow$ FFT $\Rightarrow F_0[k]$
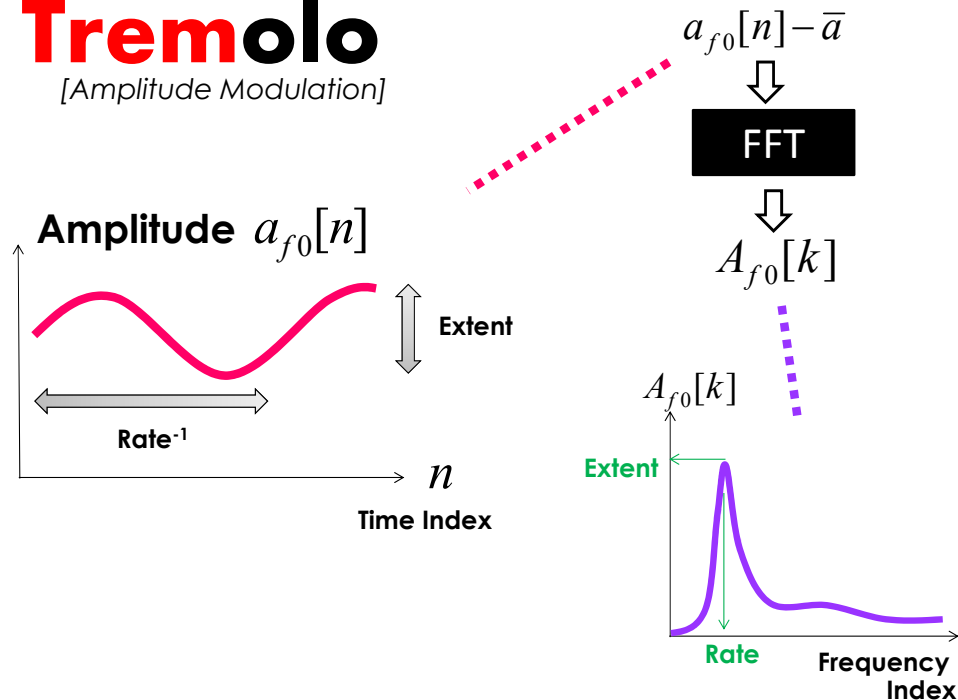
# **Vibrat**o
*[Frequency Modulation]*

# **Trem**olo
*[Amplitude Modulation]*

- Variation of sound intensity.
- **Extent of vibrato**
   = Amount of intensity variation
- **Rate of vibrato**
   = Speed which the intensity is varied
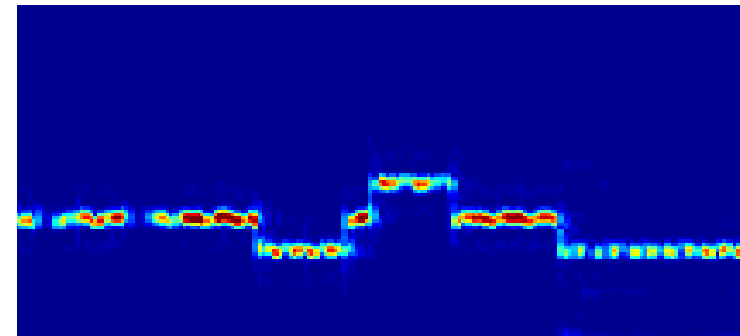- Wind and brass instruments  produce AM dominant sounds.

# **Trem**olo
*[Amplitude Modulation]*

**Amplitude** $a_{f0}[n]$

Extent

Rate$^{-1}$

$n$

**Time Index**

$a_{f0}[n] - \bar{a}$

⇩

**FFT**

⇩

$A_{f0}[k]$

$A_{f0}[k]$

Extent

Rate

**Frequency Index**

# **Trem**olo
*[Amplitude Modulation]*

# Linear Predictive Coding

• A model used for predicting a value of current sample of signal from the previous samples

• Used in Speech Analysis, Speech Synthesis, Audio Classification, Audio Compression

**Current sample**     **Previous samples**

$$s[n] = \sum_{k=1}^{P} a_k s[n-k] + \varepsilon[n]$$

**LPC Coefficients**     **Prediction Error**

---

# Linear Predictive Coding

$s[n]$

Training Signal

$\{a_k\}$

**Training LPC**

• Find $\{a_k\}$

• To minimize $\sum_n \varepsilon^2[n]$

$$\varepsilon[n] = s[n] - \hat{s}[n]$$

$$\hat{s}[n] = \sum_{k=1}^{P} a_k s[n-k]$$

---

# Linear Predictive Coding

$$E = \sum_n \varepsilon^2[n] = \sum_n \left( s[n] - \sum_{k=1}^{P} a_k s[n-k] \right)^2$$

$$\frac{\partial E}{\partial a_i} = \sum_n \left( -2s[n-i] \right) \left( s[n] - \sum_{k=1}^{P} a_k s[n-k] \right)$$

$$0 = \sum_n s[n]s[n-i] - \sum_n s[n-i] \sum_{k=1}^{P} a_k s[n-k]$$

$$0 = \sum_n s[n]s[n-i] - \sum_{k=1}^{P} a_k \sum_n s[n-k]s[n-i]$$

---

# Linear Predictive Coding

$$0 = \sum_n s[n]s[n-i] - \sum_{k=1}^{P} a_k \sum_n s[n-k]s[n-i]$$

$$0 = R[i] - \sum_{k=1}^{P} a_k R[k-i]$$

$$\sum_{k=1}^{P} a_k R[k-i] = R[i]$$

Autocorrelation Function    $R[i-j] = R[j-i] = \dfrac{\sum_n s[n-i]s[n-j]}{\sum_n s^2[n]}$

# Linear Predictive Coding

$$\frac{\partial E}{\partial a_1} = 0 \Rightarrow \sum_{k=1}^{P} a_k R[k-1] = R[1]$$

$$\frac{\partial E}{\partial a_2} = 0 \Rightarrow \sum_{k=1}^{P} a_k R[k-2] = R[2]$$

$$\vdots$$

$$\frac{\partial E}{\partial a_P} = 0 \Rightarrow \sum_{k=1}^{P} a_k R[k-P] = R[P]$$
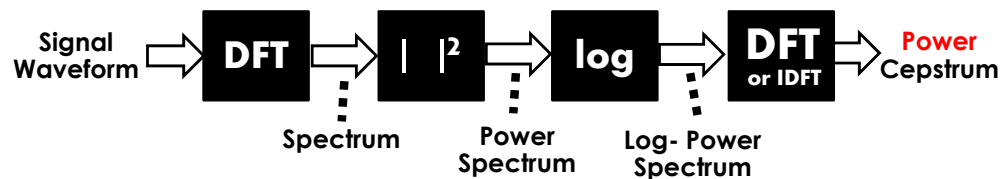
# Linear Predictive Coding

$$\begin{bmatrix} R[0] & R[1] & R[2] & \cdots & R[P-1] \\ R[1] & R[0] & R[1] & \cdots & R[P-2] \\ R[2] & R[1] & R[0] & \cdots & R[P-3] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ R[P-1] & R[P-2] & R[P-3] & \cdots & R[0] \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_P \end{bmatrix} = \begin{bmatrix} R[1] \\ R[2] \\ R[3] \\ \vdots \\ R[P] \end{bmatrix}$$

$$\downarrow$$

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_P \end{bmatrix} = \begin{bmatrix} R[0] & R[1] & R[2] & \cdots & R[P-1] \\ R[1] & R[0] & R[1] & \cdots & R[P-2] \\ R[2] & R[1] & R[0] & \cdots & R[P-3] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ R[P-1] & R[P-2] & R[P-3] & \cdots & R[0] \end{bmatrix}^{-1} \begin{bmatrix} R[1] \\ R[2] \\ R[3] \\ \vdots \\ R[P] \end{bmatrix}$$

# Cepstrum

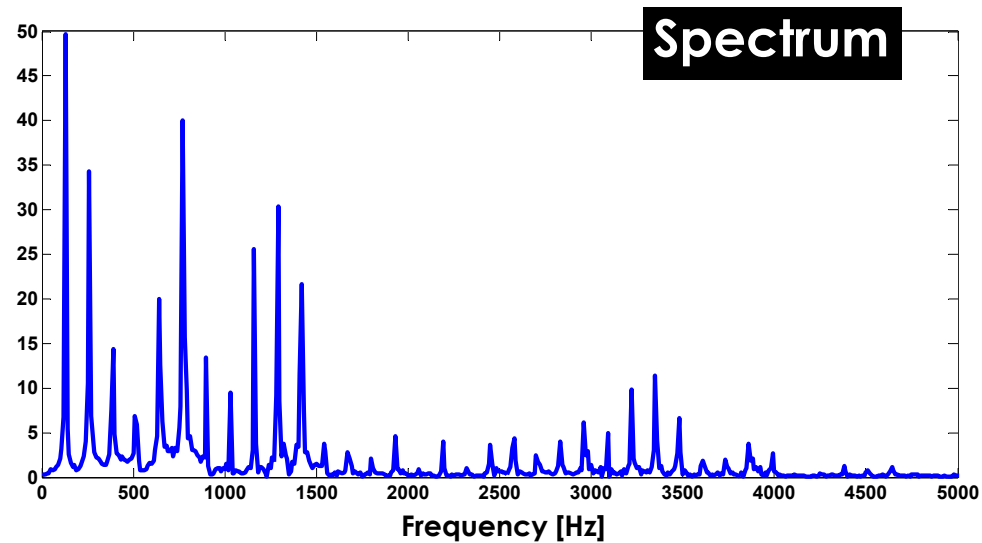**Fourier transform of the logarithm of the spectrum of signal**

Signal Waveform → **DFT** → **| |²** → **log** → **DFT or IDFT** → Power Cepstrum

Spectrum    Power Spectrum    Log- Power Spectrum

Spectrum ⇒ Cepstrum

Frequency ⇒ Quefrency

Filter ⇒ Lifter

# Cepstrum



Waveform a:
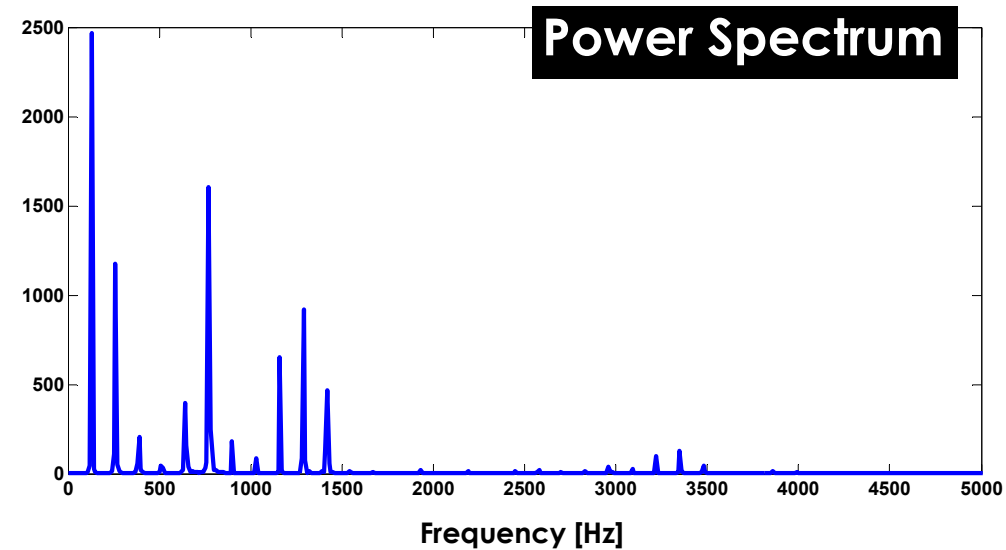
Time [sec]

# Cepstrum

Spectrum

# Cepstrum

Power Spectrum

# Cepstrum

Log Power Spectrum
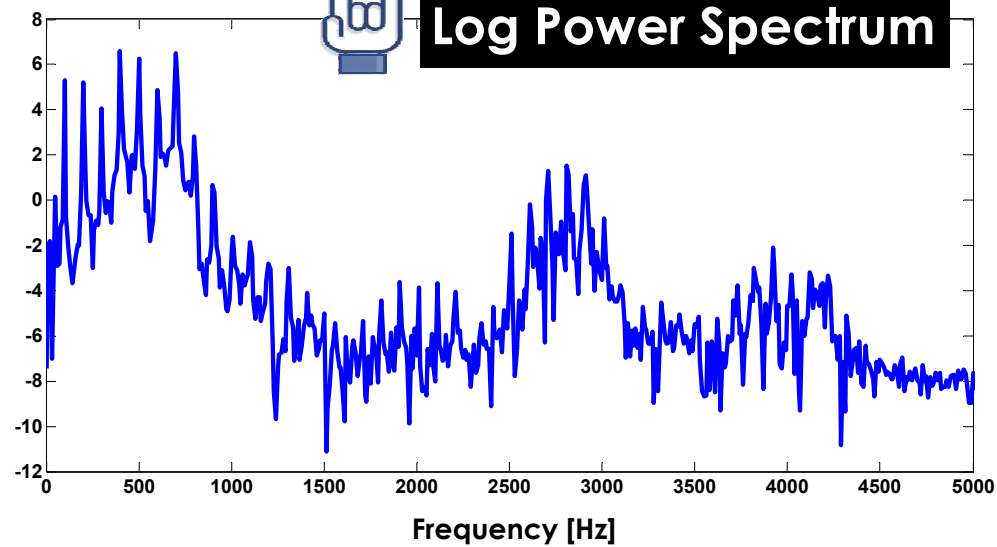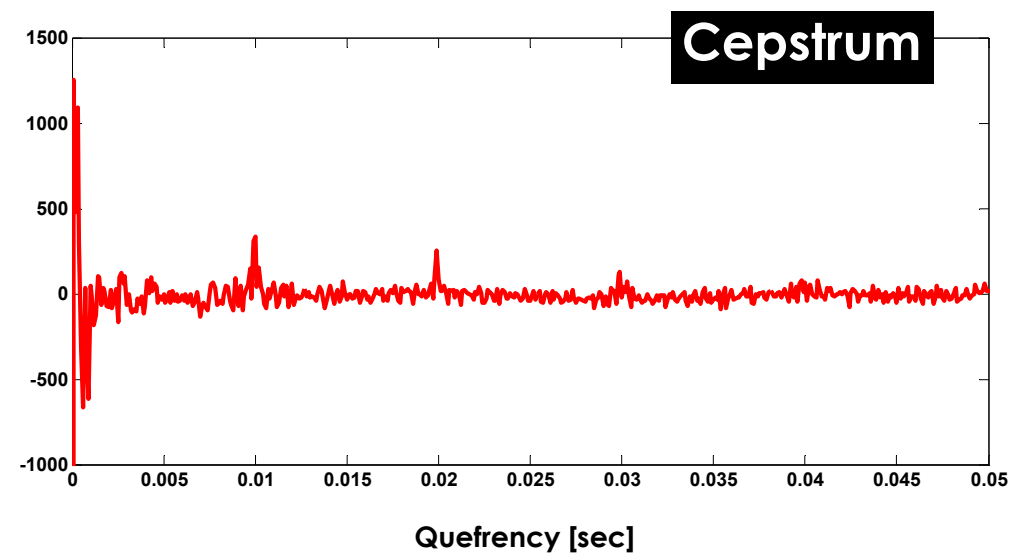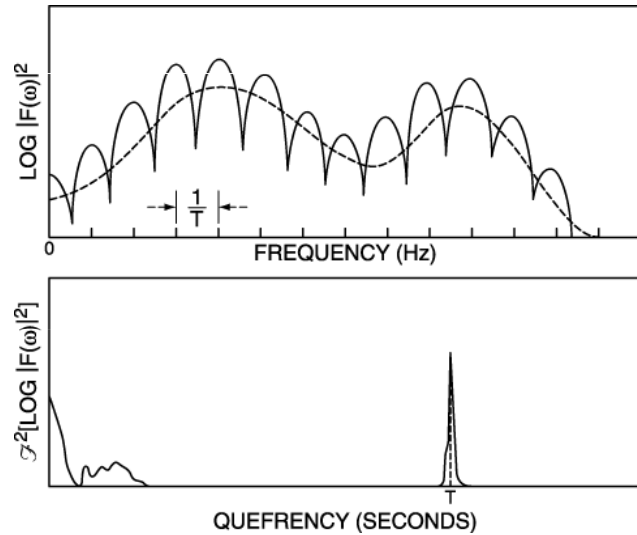
# Cepstrum

Cepstrum

# Cepstrum

**Waveform o:**



Time [sec]

# Cepstrum

**Spectrum**



Frequency [Hz]

# Cepstrum

**Log Power Spectrum**



Frequency [Hz]

# Cepstrum

**Cepstrum**



Quefrency [sec]

# Cepstrum

LOG $|F(\omega)|^2$

$\leftarrow \frac{1}{T} \rightarrow$

0   FREQUENCY (Hz)
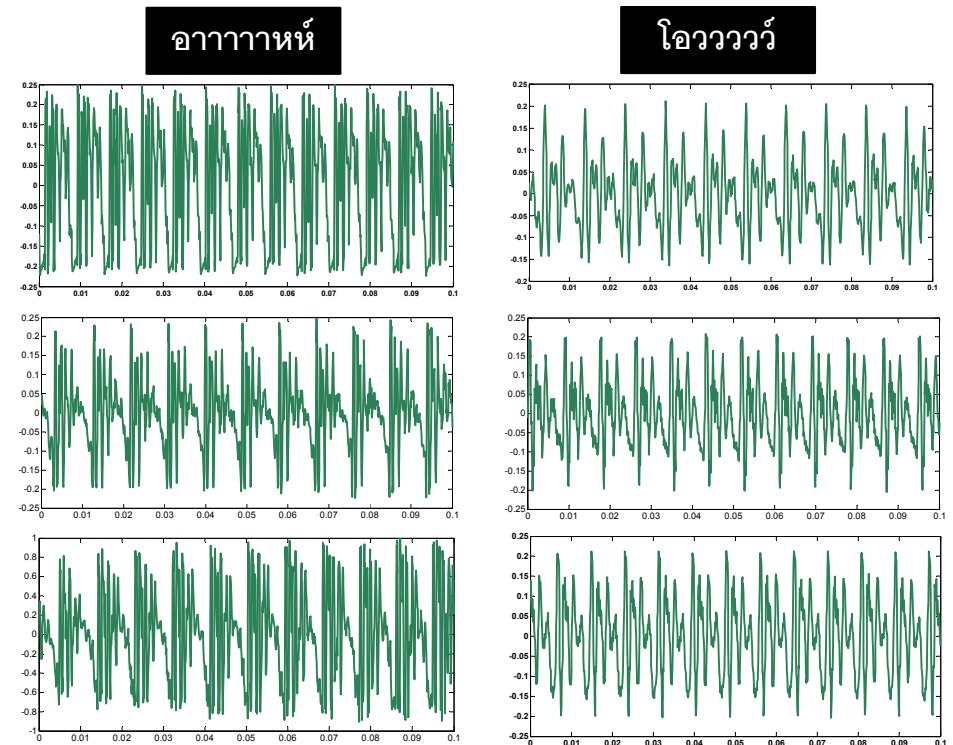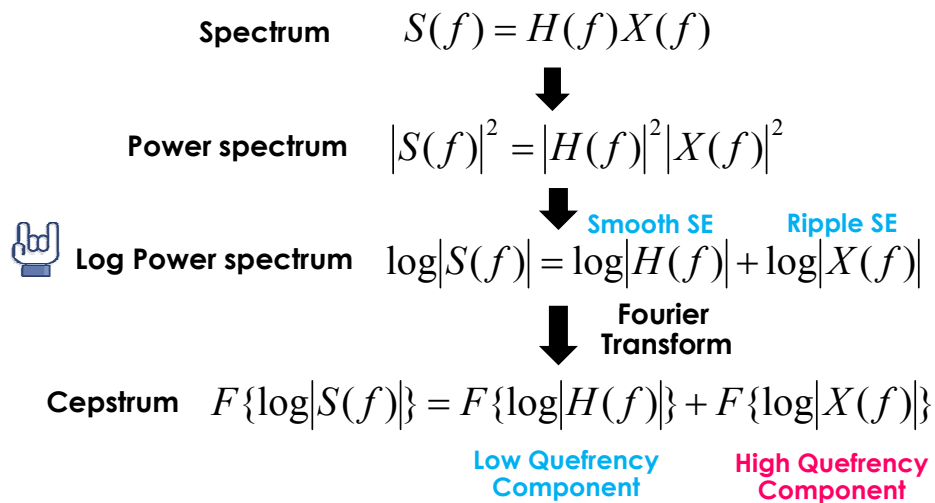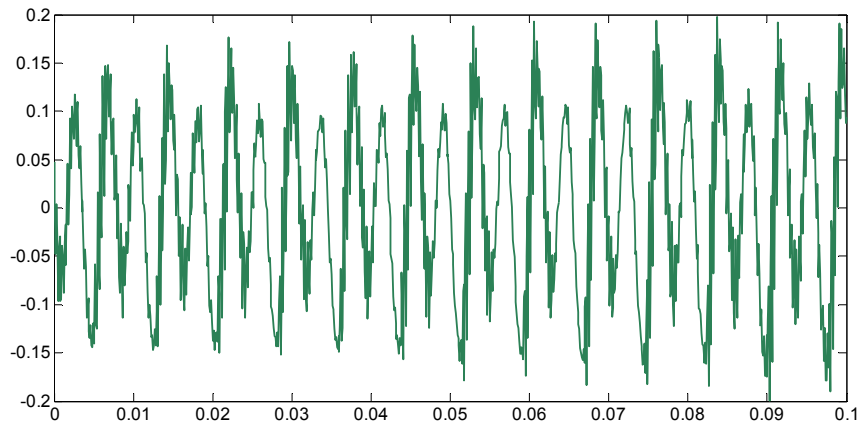
$\mathcal{F}^2[$LOG $|F(\omega)|^2]$

T

QUEFRENCY (SECONDS)

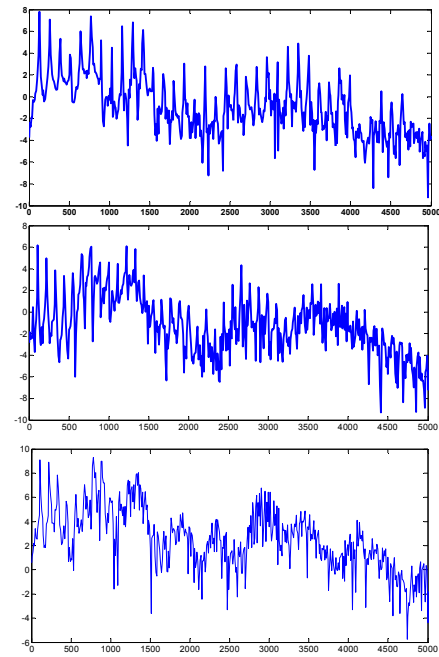http://postgavindisorder.blogspot.com/2010_09_01_archive.html

# Cepstrum

$x(t)$ ➡ **System** ➡ $s(t)$

Impulse Response $h(t)$

$$s(t) = h(t) * x(t)$$

**Fourier Transform** ⬇

$$S(f) = H(f)X(f)$$

$X(f)$

**Spectral envelope = Ripple**

$f$

$H(f)$

**Spectral envelope = Smooth**

$f$

$S(f)$

$f$

# Cepstrum

Spectrum $\qquad S(f) = H(f)X(f)$

⬇

Power spectrum $\qquad |S(f)|^2 = |H(f)|^2 |X(f)|^2$

⬇ Smooth SE    Ripple SE

Log Power spectrum $\qquad \log|S(f)| = \log|H(f)| + \log|X(f)|$

⬇ **Fourier Transform**

Cepstrum $\qquad F\{\log|S(f)|\} = F\{\log|H(f)|\} + F\{\log|X(f)|\}$

**Low Quefrency Component**     **High Quefrency Component**

อาาาาห์     โอวววว์

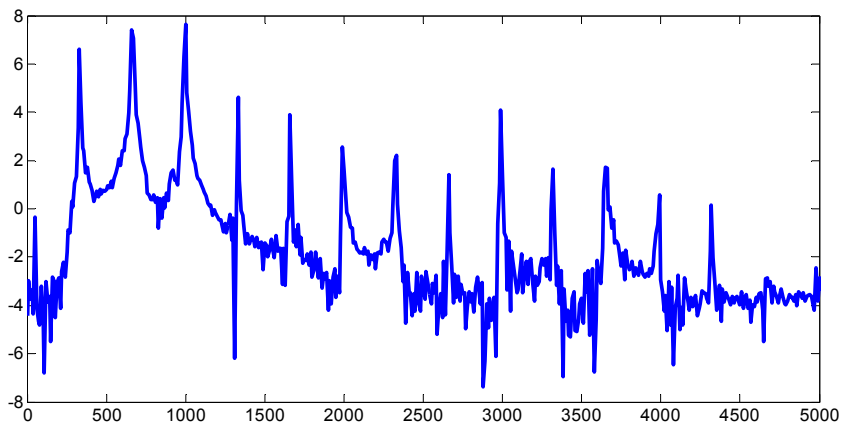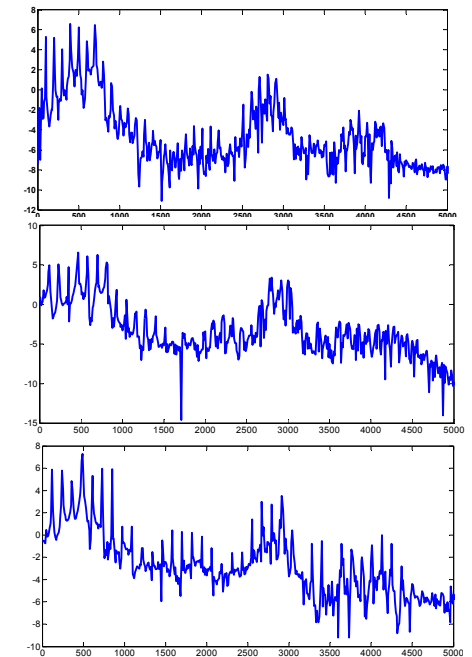ให้ทายว่าเสียงอะไร?

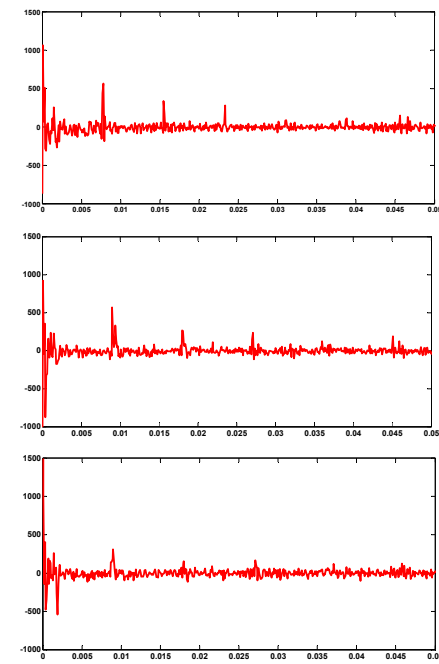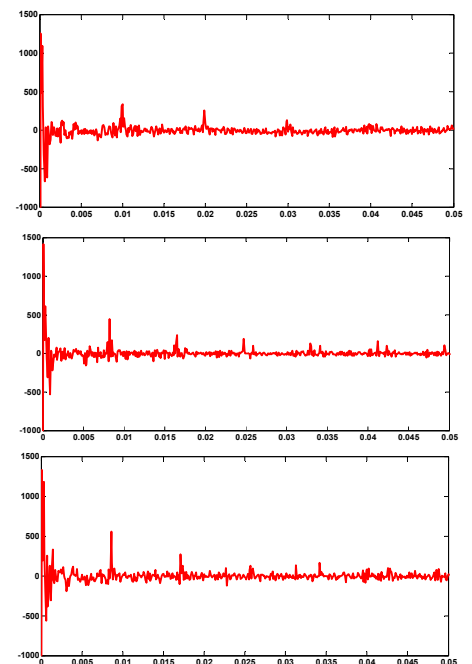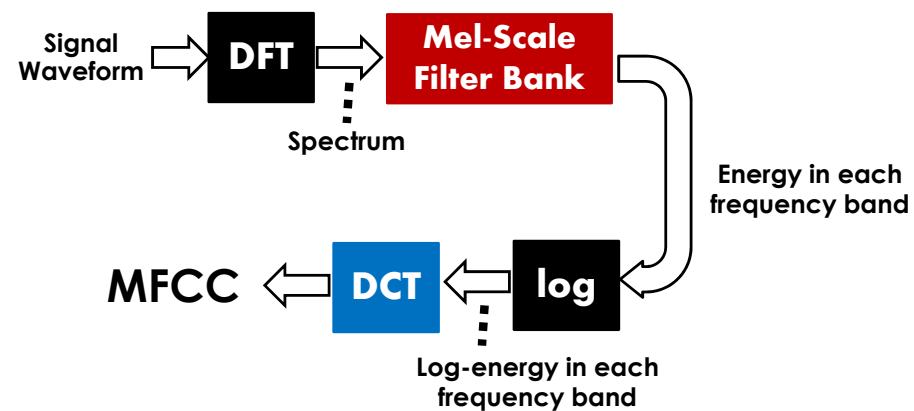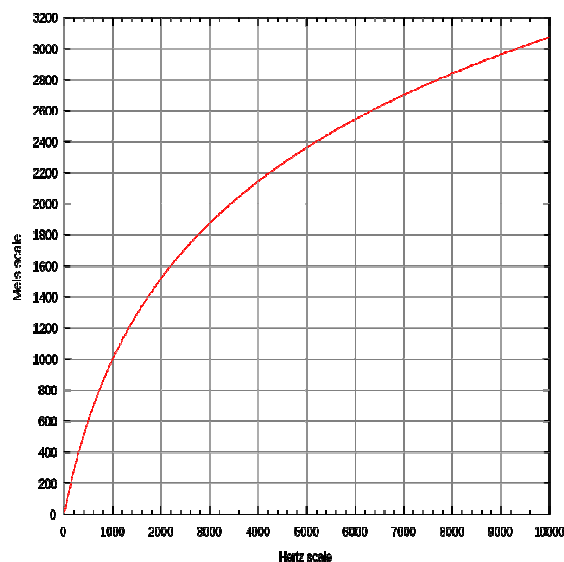อาาาาหน์   โอวววว์

ให้ทายว่าเสียงอะไร?

อาาาาหน์   โอวววว์

**ให้ทายว่าเสียงอะไร?**

# Mel-Frequency Cepstrum Coefficient

## [MFCC]



Signal Waveform → DFT → Mel-Scale Filter Bank

Spectrum

Energy in each frequency band

MFCC ← DCT ← log

Log-energy in each frequency band

**DCT = Discrete Cosine Transform**

# Mel-Scale



$$f_{[mel]} = 1127.0148 \ln\left(1 + \frac{f_{[Hz]}}{700}\right)$$

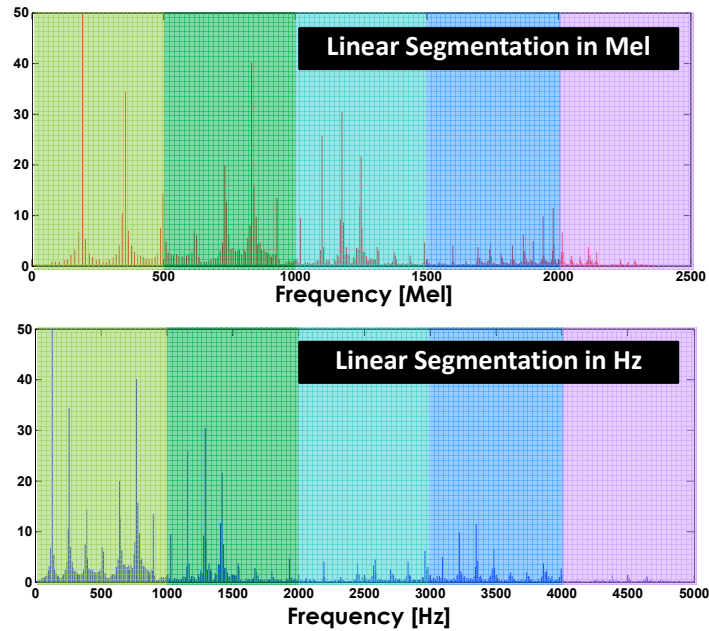• **Perceptual scale of pitches**

| Hertz | Mel |
|-------|------|
| 174 | 250 |
| 391 | 500 |
| 662 | 750 |
| 1000 | 1000 |
| 1949 | 1500 |
| 3429 | 2000 |
| 5734 | 2500 |

http://en.wikipedia.org/wiki/File:Mel-Hz_plot.svg

# Mel-Scale

# Mel-Scale



Linear Segmentation in Mel

Linear Segmentation in Hz

# Mel-Scale



Linear Segmentation in Mel

Non-linear Segmentation in Hz
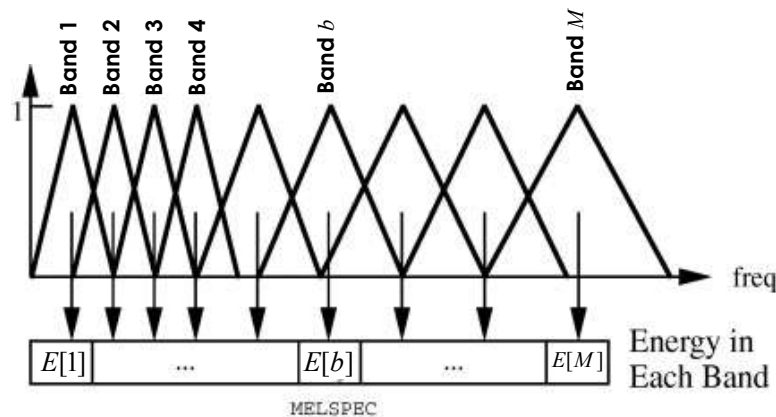
# Mel-Scale
# Filter Bank

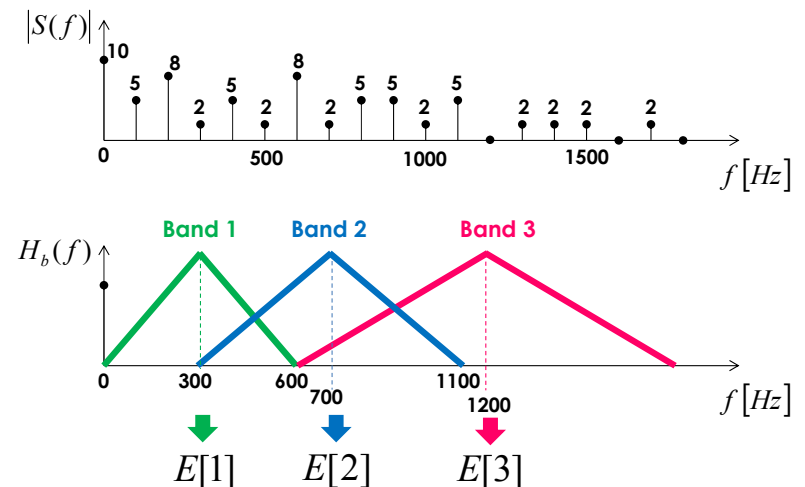$$E[b] = \sum_{f \in Band[b]} |S(f)| H_b(f)$$

$E[b]$ : Energy in Band $b$
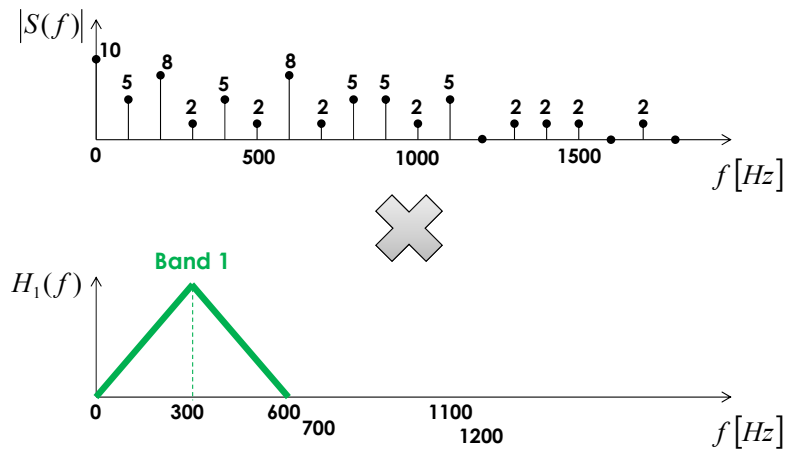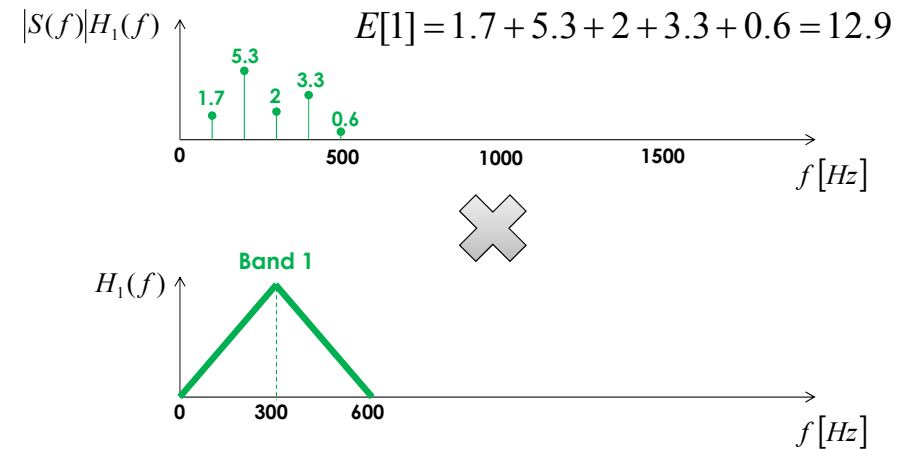
$H_b(f)$ : Freq. Response of Filter for Band $b$



Energy in Each Band

MELSPEC

Unsupervised speaker segmentation with residual phase and MFCC features
S. Jothilakshmi, , V. Ramalingam , S. Palanivel

# Mel-Scale
# Filter Bank

# **Mel**-Scale Filter Bank

$$E[1] = \sum_{f \in Band[1]} |S(f)| H_1(f)$$

# **Mel**-Scale Filter Bank

$$E[1] = \sum_{f \in Band[1]} |S(f)| H_1(f)$$

$$E[1] = 1.7 + 5.3 + 2 + 3.3 + 0.6 = 12.9$$

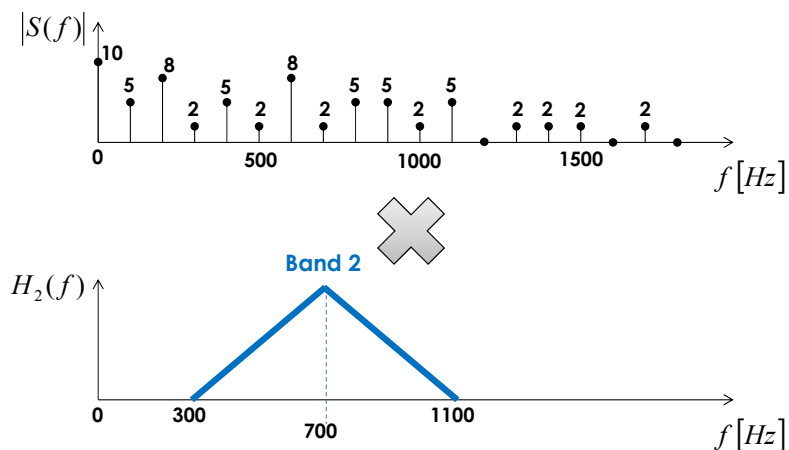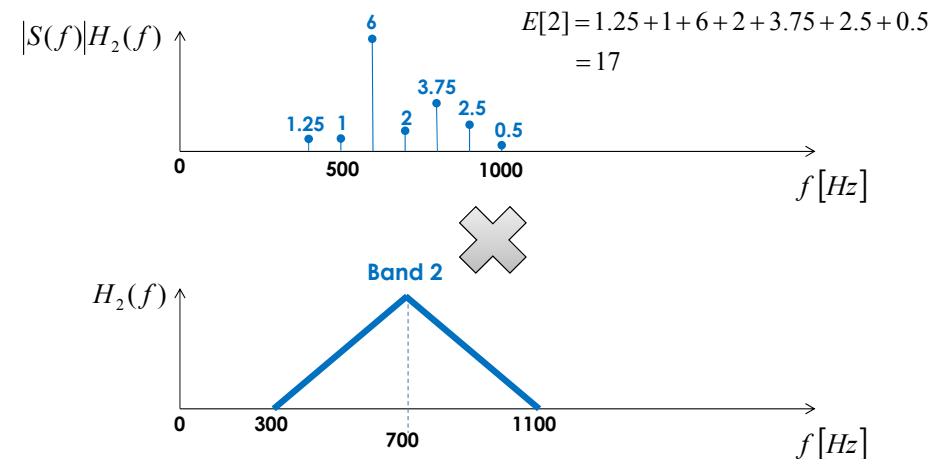# **Mel**-Scale Filter Bank

$$E[2] = \sum_{f \in Band[2]} |S(f)| H_2(f)$$

# **Mel**-Scale Filter Bank

$$E[2] = \sum_{f \in Band[2]} |S(f)| H_2(f)$$

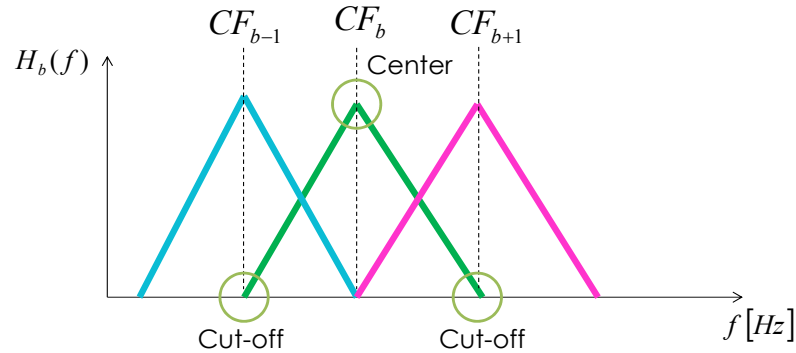$$E[2] = 1.25 + 1 + 6 + 2 + 3.75 + 2.5 + 0.5 = 17$$

# **Mel**-Scale
# **Filter Bank**

- Cut-off frequency of filter in the current band determined by the center frequencies of the two adjacent filters.



$CF_{b-1}$  $CF_b$  $CF_{b+1}$

$H_b(f)$

Center

Cut-off  Cut-off

$f[Hz]$

---

# **Mel**-Scale
# **Filter Bank**

- Triangular Filter

$$H_b(f) = \begin{cases} \dfrac{f - CF_{b-1}}{CF_b - CF_{b-1}}; & CF_{b-1} < f < CF_b \\[2mm] \dfrac{f - CF_{b+1}}{CF_b - CF_{b+1}}; & CF_b \le f < CF_{b+1} \\[2mm] 0; & ; otherwise \end{cases}$$

$CF_b = $ Center Frequency of Band $b[Hz]$

---

# **Mel**-Scale
# **Filter Bank**

- Center Frequency

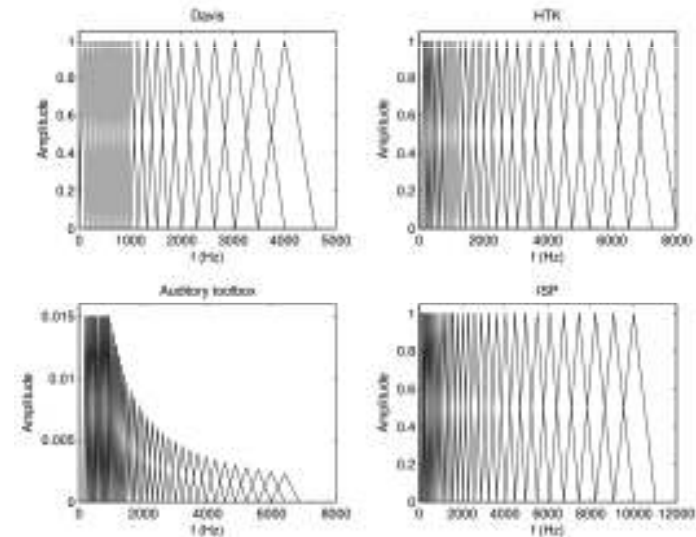$$CF_b = F_{min} + b \cdot \frac{F_{max} - F_{min}}{M+1} \qquad in \; [mel]$$

⬇

$$CF_b \; in \; [Hz]$$

$M = $ Number of Bands

$F_{min} \sim F_{max} = $ Frequency Range

---

# **Mel**-Scale
# **Filter Bank**

# **D**iscrete **C**osine **T**ransform
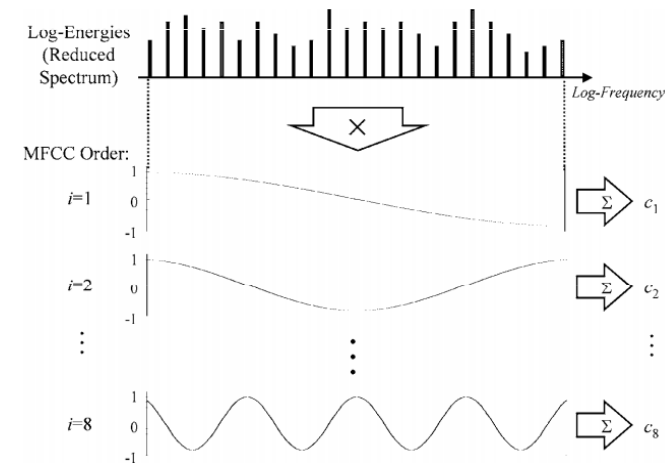## [**DCT**]

$$c[i] \; = \sum_{b=1}^{M} \ln(E[b])\cos\left( i\left( b - \frac{1}{2} \right)\frac{\pi}{M} \right)$$

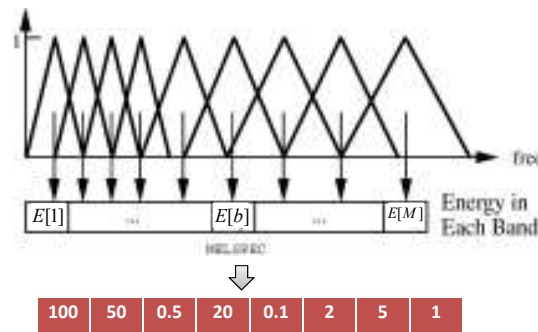$c[i]$ = M*e*l Frequency Cepstrum Coefficient

$i$ = Order of MFCC

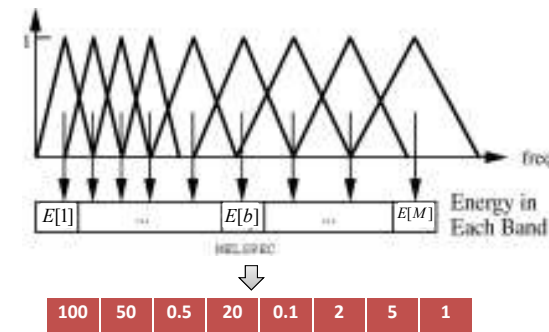$M$ = Number of Frequency Bands

# **D**iscrete **C**osine **T**ransform
## [**DCT**]



Log-Energies (Reduced Spectrum)

Log-Frequency

MFCC Order:

$i$=1 $\quad \Sigma \quad c_1$

$i$=2 $\quad \Sigma \quad c_2$

$i$=8 $\quad \Sigma \quad c_8$

# **D**iscrete **C**osine **T**ransform
## [**DCT**]



$E[1]$ ... $E[b]$ ... $E[M]$ Energy in Each Band

| 100 | 50 | 0.5 | 20 | 0.1 | 2 | 5 | 1 |

$$c[1] \; = \sum_{b=1}^{8} \ln(E[b])\cos\left( \left( b - \frac{1}{2} \right)\frac{\pi}{8} \right)$$

$$= \ln(100)\cos\left(\frac{\pi}{16}\right) + \ln(50)\cos\left(\frac{3\pi}{16}\right) + \ln(0.5)\cos\left(\frac{5\pi}{16}\right) + ... + \ln(1)\cos\left(\frac{15\pi}{16}\right)$$

$$= 6.6947$$

# **D**iscrete **C**osine **T**ransform
## [**DCT**]



$E[1]$ ... $E[b]$ ... $E[M]$ Energy in Each Band

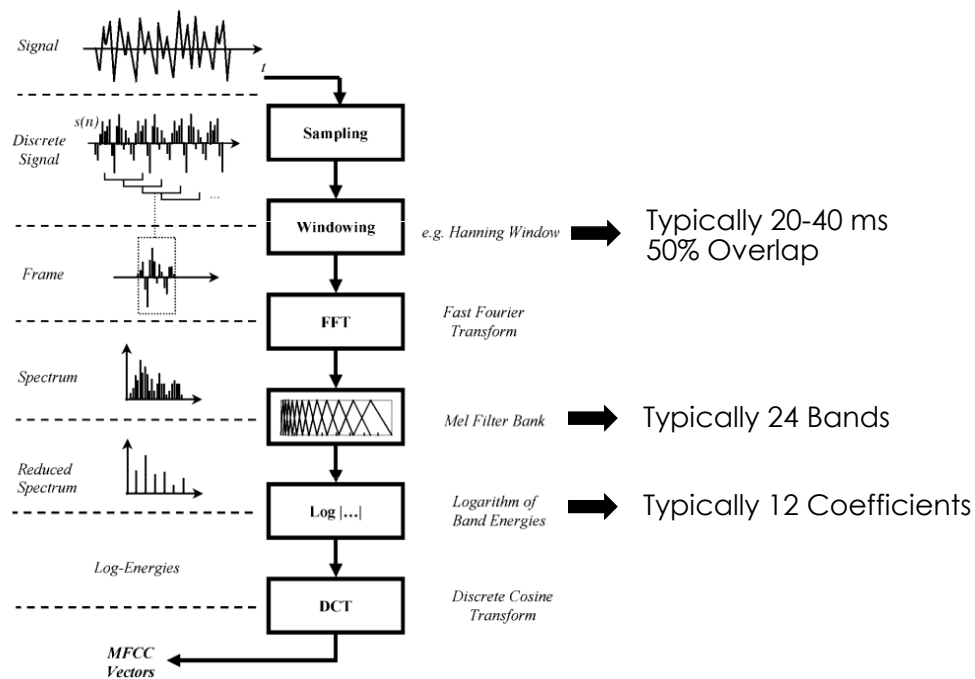| 100 | 50 | 0.5 | 20 | 0.1 | 2 | 5 | 1 |

$$c[2] \; = \sum_{b=1}^{8} \ln(E[b])\cos\left( 2\left( b - \frac{1}{2} \right)\frac{\pi}{8} \right)$$

$$= \ln(100)\cos\left(\frac{\pi}{8}\right) + \ln(50)\cos\left(\frac{3\pi}{8}\right) + \ln(0.5)\cos\left(\frac{5\pi}{8}\right) + ... + \ln(1)\cos\left(\frac{15\pi}{8}\right)$$
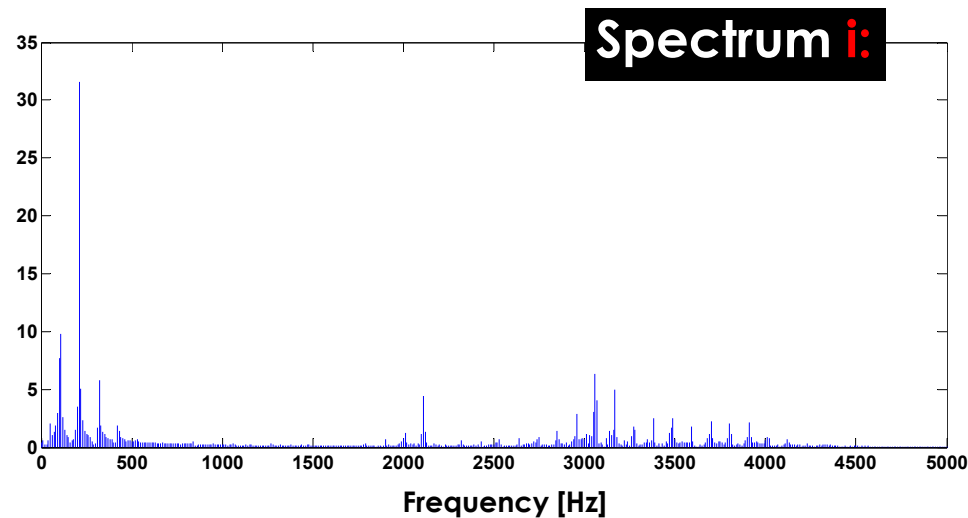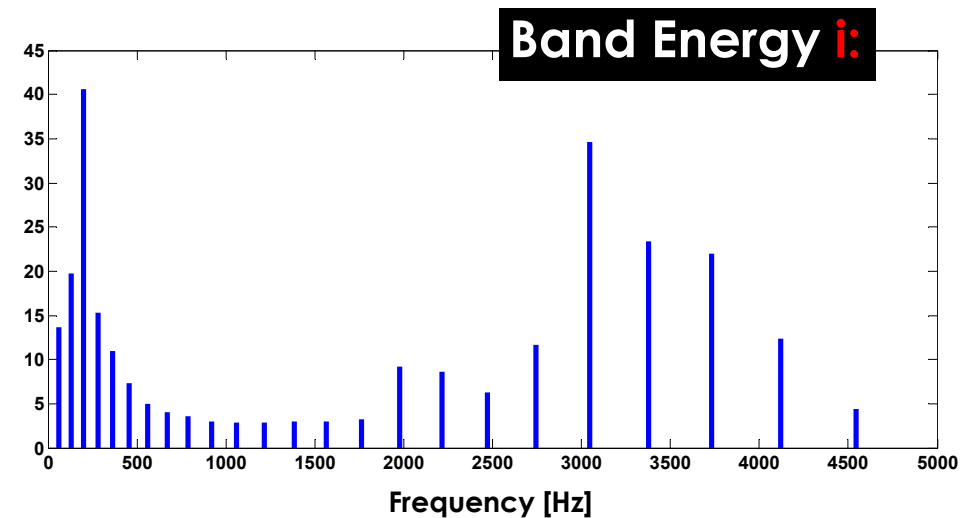
$$= 5.7272$$

Signal

Discrete Signal $s(n)$

Sampling

Frame

Windowing    e.g. Hanning Window → **Typically 20-40 ms 50% Overlap**

FFT    Fast Fourier Transform

Spectrum

Mel Filter Bank → **Typically 24 Bands**

Reduced Spectrum

Log |...|    Logarithm of Band Energies → **Typically 12 Coefficients**

Log-Energies

DCT    Discrete Cosine Transform

MFCC Vectors

MPEG-7 Audio and. Beyond. Audio Content Indexing and. Retrieval. Hyoung-Gook Kim, et. al.

# **Mel-Frequency Cepstrum Coefficient**

**Waveform i:**



Time [sec]

# **Mel-Frequency Cepstrum Coefficient**

**Spectrum i:**



Frequency [Hz]

# **Mel-Frequency Cepstrum Coefficient**

**Band Energy i:**



Frequency [Hz]

# Mel-Frequency Cepstrum Coefficient

Band Log-Energy i:

# Mel-Frequency Cepstrum Coefficient

MFCC i:

อี

อู

อี

อู

อี

อู

อี

อู

โอ

อา

แอ

เอ

กากจุงเบย

กากขิง ๆ

มันกากมาก

**Time Frame**

หล่อ

พ่องง