

# I. Excitation Signal Modeling

$$e[t] = a[t] * x[t]$$

↳ glottal pulse :



Pitch

$$T_0 = \frac{1}{\text{Pitch}}$$

LPC analysis of residual signal is glottal pulse is represented

① pitch extraction

pitch extraction is correlation of  $\frac{1}{2}$  sec

$$\text{auto correlation } ac(\tau) = \sum_{n=0}^{N-\tau} x[n] \cdot x[n+\tau]$$

\* auto correlation of  $\frac{1}{2}$  sec.  $\tau = 0.5$  sec.  $\tau = 0.5$  sec.  $\tau = 0.5$  sec.  $\tau = 0.5$  sec.

ac ← auto correlation of  $\frac{1}{2}$  sec. peak is max of pitch is  $\frac{1}{2}$  sec.

(voiced) residual of  $x(t)$  preemphasis (x)  $x(t)$  (0)



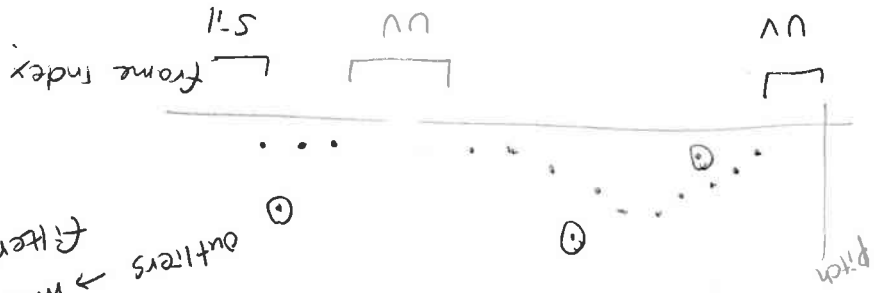
(unvoiced, H1 (Gaussian random)) glottal pulse is pitch of  $\frac{1}{2}$  sec

⊗ V/UV/silence detection of  $\frac{1}{2}$  sec

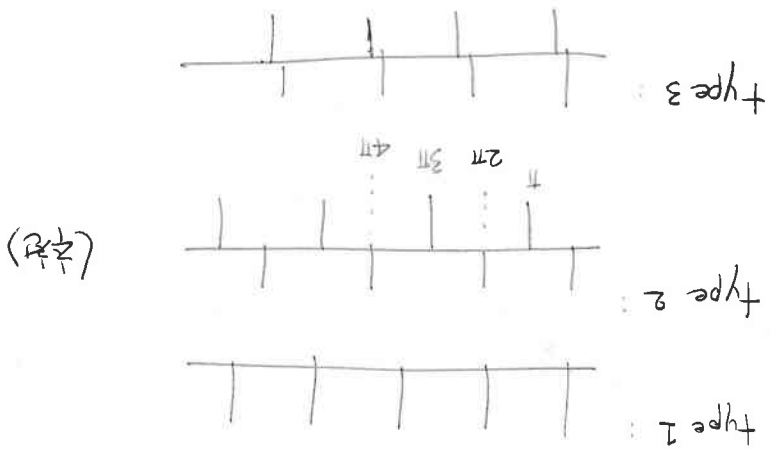
a)  $\frac{ac[\tau]}{ac[0]} \geq \theta$  → V  
 $\frac{ac[\tau]}{ac[0]} < \theta$  → UV  
 decision rule: voiced is vocal cord is vibrating high frequency energy in  $\frac{1}{2}$  sec. UV is aperiodic and  $\frac{1}{2}$  sec. no preemphasis.

b) project 2014  $\frac{1}{2}$  sec V/UV/s classifier After

① outliers → median filter  $\frac{5 \times 5}{25}$  사용

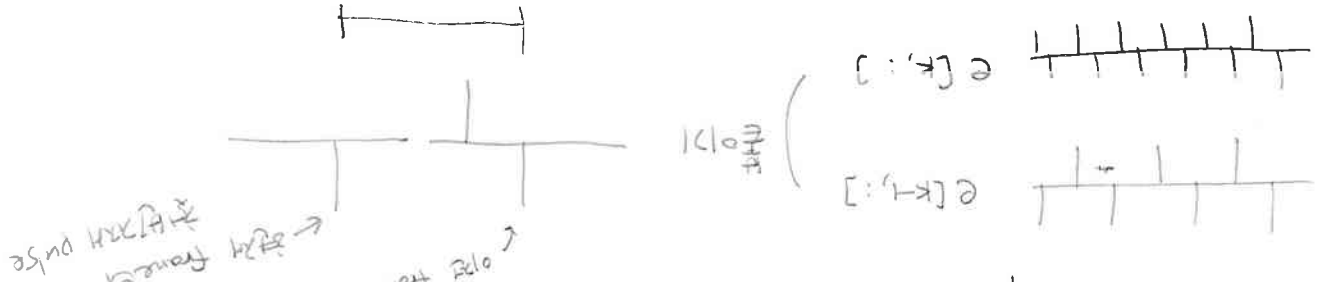


② global pulse AHJ



이 pulse global pulse에 대해 각각 다르게 처리

③ Concatenation to previous excitation : 마지막 observation



3-2 Concatenation

은  $\frac{1}{2}\pi$ 가 주파수

이 pulse residual과  $\pi$ 가  $\pi$ 와 align되는 위치를 찾는 방법

$$r[k] = a[k] * x[k] \quad (\text{residual})$$

$$e(t, z) = \frac{z}{z^2} = z = 0 \sim \frac{z}{p}$$

$$z^* = \arg \max_z e(t, z) \cdot r[k]$$

cross-correlation with residual

$$\arg \max_z \text{sign}(e(t, z)) \cdot \text{sign}(r[k]) \quad \frac{1}{2}\pi \text{의 cross-correlation}$$

Adapted pulse energy residual signal at  $2\pi L \frac{t}{T}$  scaling

$$g[k] = \sqrt{\frac{1}{N_p} \sum_{n=0}^{N_p-1} r^2(t)} \quad (\text{residual standard deviation})$$

$$\hat{e}[k, t] \leftarrow g[k] \frac{\sqrt{\frac{1}{N_p} \sum_{n=0}^{N_p-1} r^2(t)}}{\sqrt{N_p \sum_{n=0}^{N_p-1} r^2(t)}} \quad (\text{frame gain } g[k] \text{ at } t \text{ is parameter } 2 \text{ at } t)$$

⑤ Unvoiced ?

$$\hat{e}[k, t] \leftarrow g[k] \cdot \text{rand. gauss.}(0,1)$$

⑥ Encoding 600 bytes/frame = 600 bytes/sec

Flagvuv byte : voiced, unvoiced, silence

$G[k]$  : byte or short - excitation gain or frame gain

$P[k]$  : short - pitch

offset[k] : byte -  $2^*$

