

# 딥러닝 기반 사용자 적응형 음성합성 시스템

VOICE & DIALOGUE (TRACK A)  
15:00 ~ 15:25

송은우

NAVER

# TTS AT NAVER

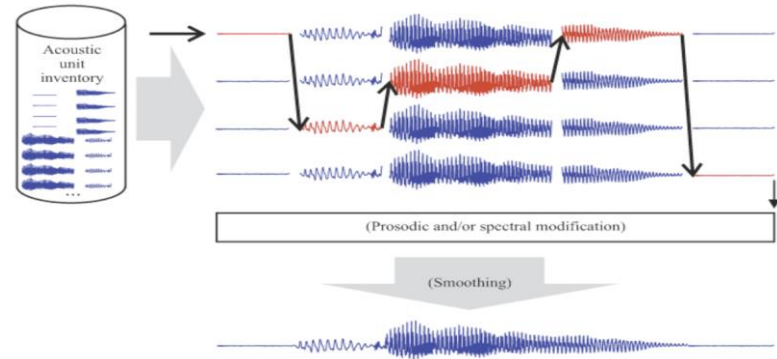
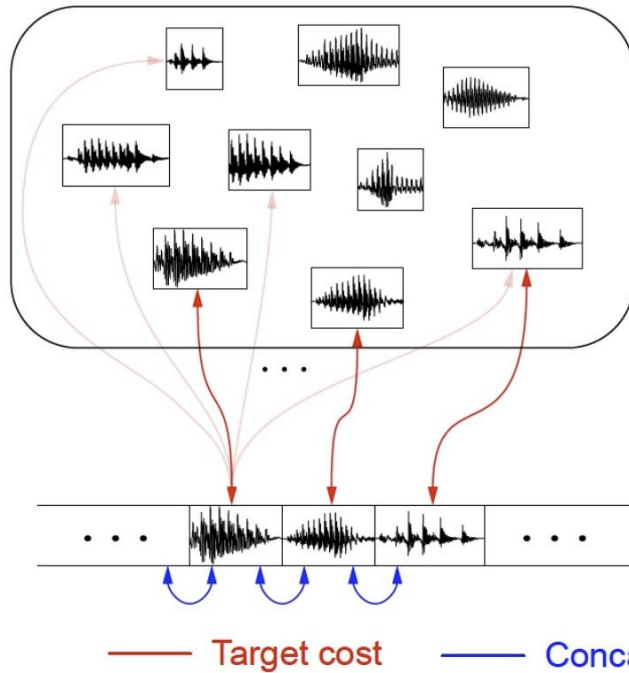
## nVoice applications

Clova	Papago	NAVER i	Whale
<p>유인나 오디오북</p> 	<p>사전 예문읽기</p> 	<p>뉴스 본문듣기</p> 	<p>Open API</p> 

# UNIT-SELECTION SPEECH SYNTHESIS SYSTEM (UTS)

Concatenate speech units (waveform) from a LARGE database

All segments



$$C(t_1^n, u_1^n) = \sum_{i=1}^n C^{target}(t_i, u_i) + \sum_{i=2}^n C^{join}(u_{i-1}, u_i)$$

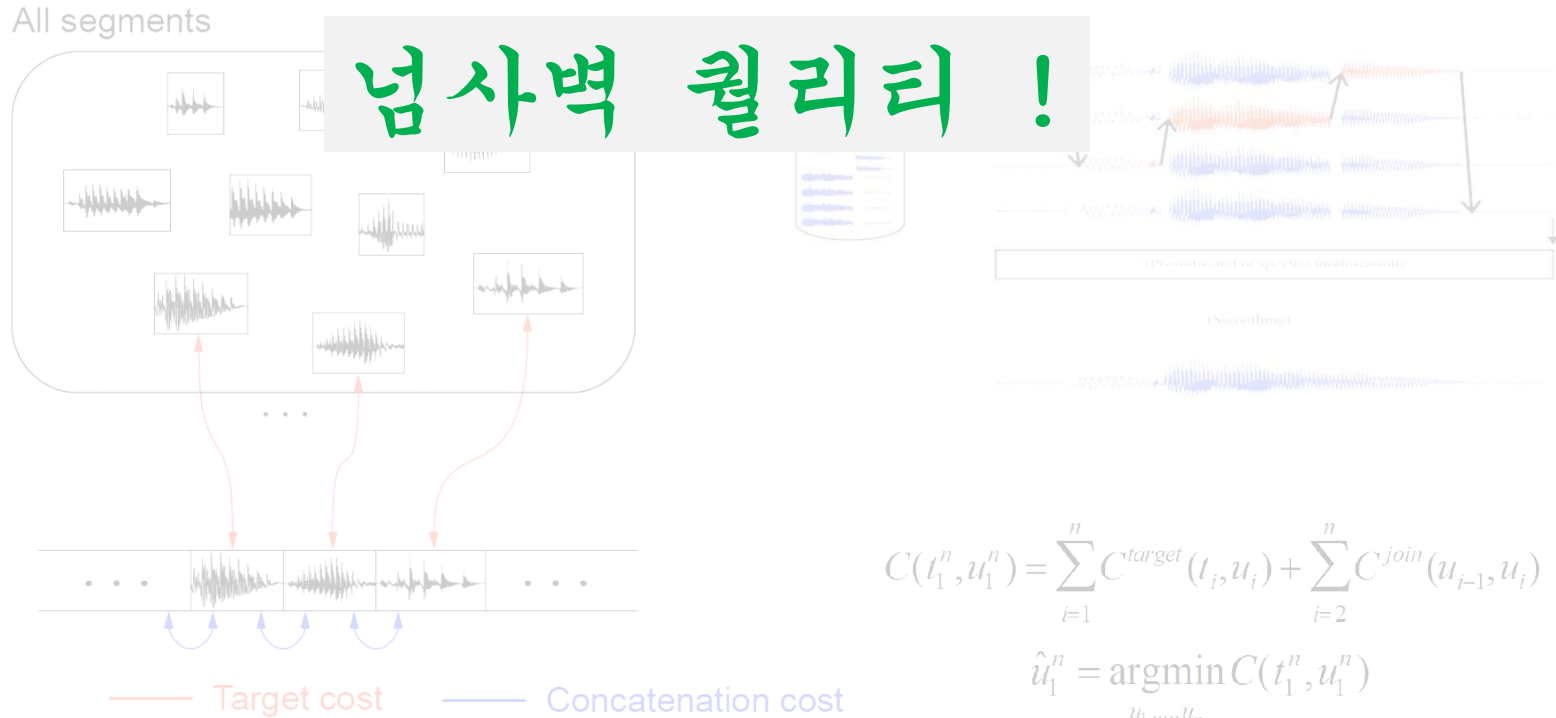
$$\hat{u}_1^n = \underset{u_1, \dots, u_n}{\operatorname{argmin}} C(t_1^n, u_1^n)$$

출처:

T. Keichi, and H. Zen. "Fundamentals and recent advances in HMM-based speech synthesis." *Tutorial of INTERSPEECH*, 2009.

# UNIT-SELECTION SPEECH SYNTHESIS SYSTEM (UTS)

Concatenate speech units (waveform) from a LARGE database

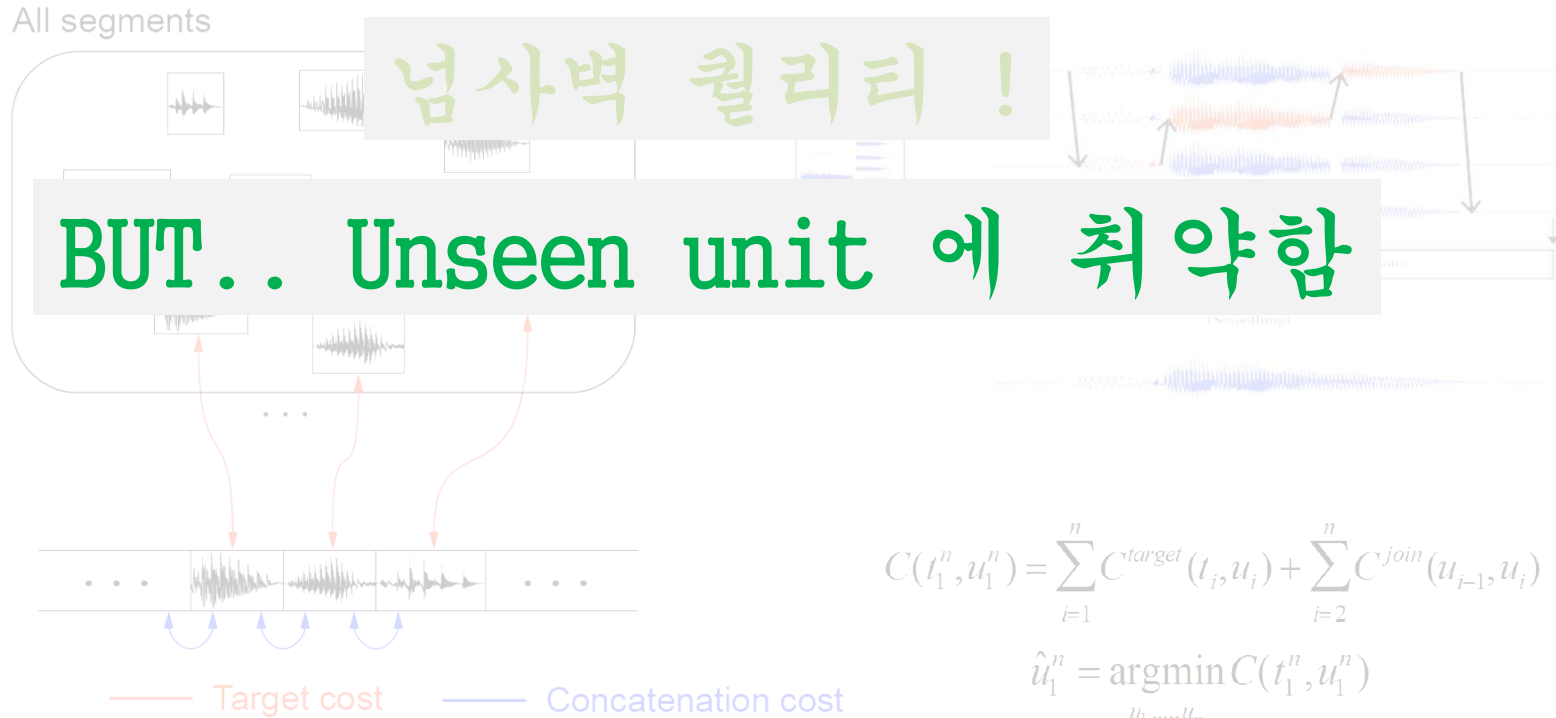


출처:

T. Keiichi, and H. Zen. "Fundamentals and recent advances in HMM-based speech synthesis." *Tutorial of INTERSPEECH*, 2009.

# UNIT-SELECTION SPEECH SYNTHESIS SYSTEM (UTS)

Concatenate speech units (waveform) from a LARGE database



출처:

T. Keiichi, and H. Zen. "Fundamentals and recent advances in HMM-based speech synthesis." *Tutorial of INTERSPEECH*, 2009.

# UNIT-SELECTION SPEECH SYNTHESIS SYSTEM (UTS)

Concatenate speech units (waveform) from a LARGE database



$$C(t_1^n, u_1^n) = \sum_{i=1}^n C^{target}(t_i, u_i) + \sum_{i=2}^n C^{join}(u_{i-1}, u_i)$$
$$\hat{u}_1^n = \underset{u_1, \dots, u_n}{\operatorname{argmin}} C(t_1^n, u_1^n)$$



출처:

T. Keiichi, and H. Zen. "Fundamentals and recent advances in HMM-based speech synthesis." *Tutorial of INTERSPEECH*, 2009.

# UNIT-SELECTION SPEECH SYNTHESIS SYSTEM (UTS)

Concatenate speech units (waveform) from a LARGE database



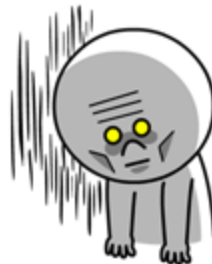
100명의 목소리를 만들고 싶다면 ?  $u_i$

— Target cost

— C

ist

$$\hat{u}_1^n = \underset{u_1, \dots, u_n}{\operatorname{argmin}} C(t_1^n, u_1^n)$$

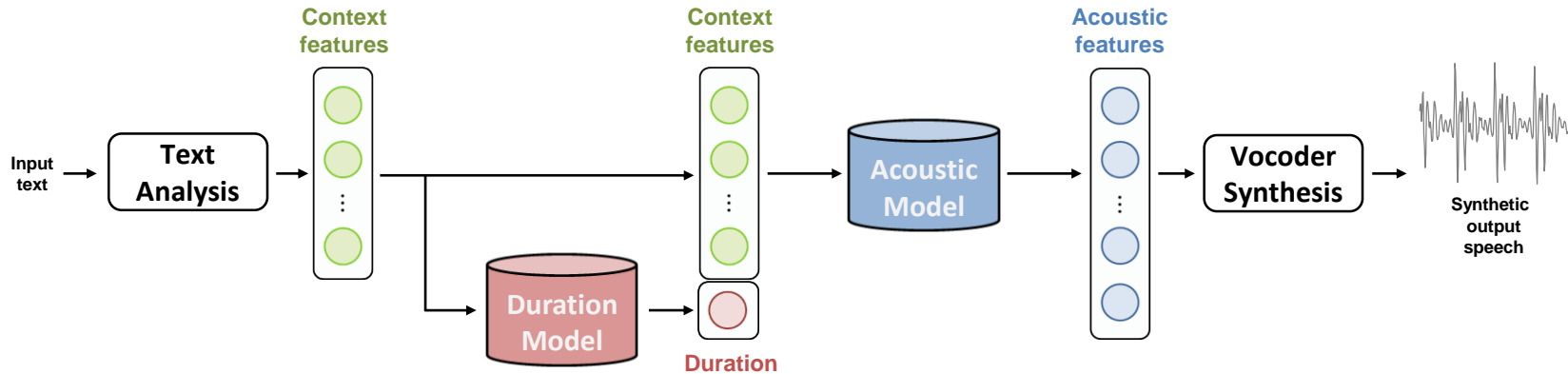


출처:

T. Keichi, and H. Zen. "Fundamentals and recent advances in HMM-based speech synthesis." *Tutorial of INTERSPEECH*, 2009.

# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Generate speech parameters from input text and deep neural network

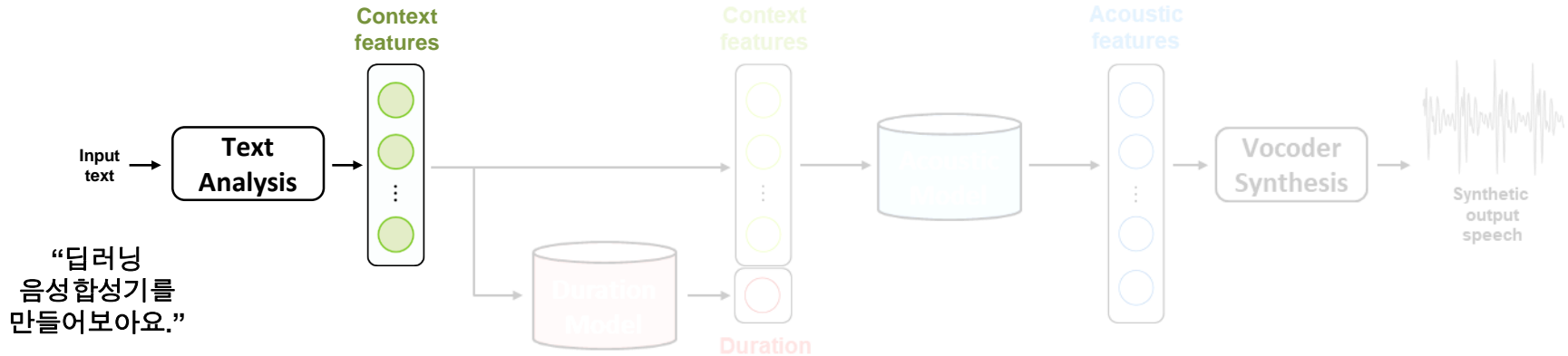


- NLP frontend : text analyzer
- Speech backend : vocoder synthesizer
- Deep learning model : duration & acoustic model



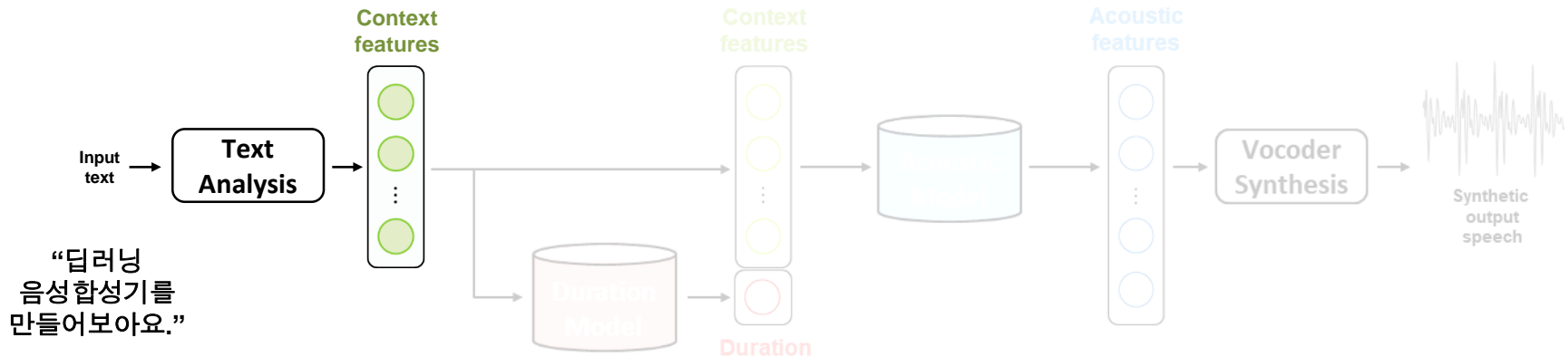
# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Text analyzer



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Text analyzer

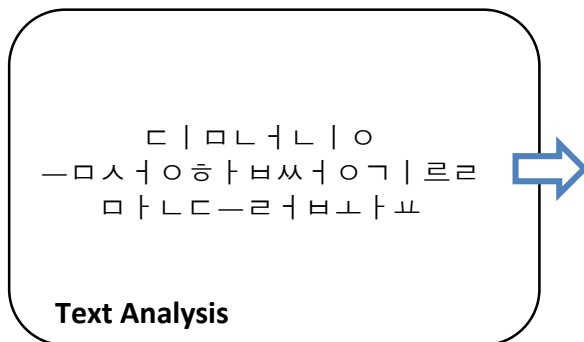
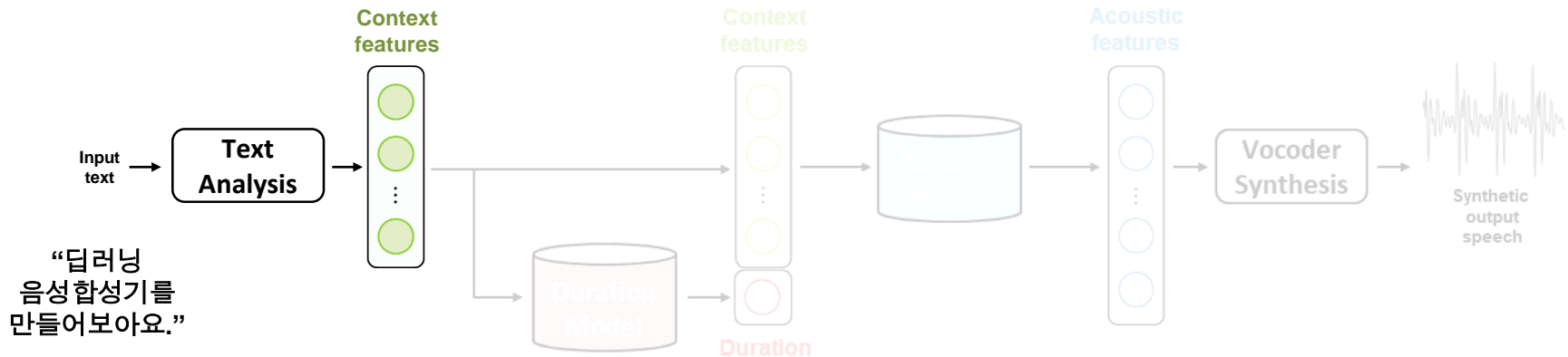


ㄷ | ㅁ | ㄴ | ㄴ | ㅇ  
- ㅁ | ㅅ | ㅇ | ㅎ | ㅈ | ㅅ | ㅇ | ㄱ | ㄹ | ㄹ  
ㅁ | ㅈ | ㄷ | - ㄹ | ㅈ | ㅅ | ㅈ | ㅈ

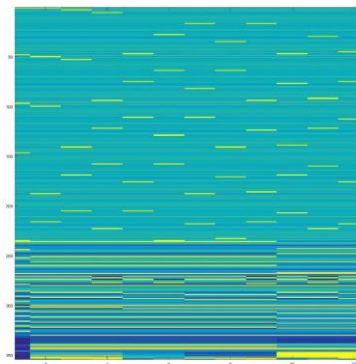
Text Analysis

# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Text analyzer



(목음) ㄷ | ㅁ | ㄴ | ㅈ | ㄴ | | ㅇ | - | ㅁ | ㅅ | ...

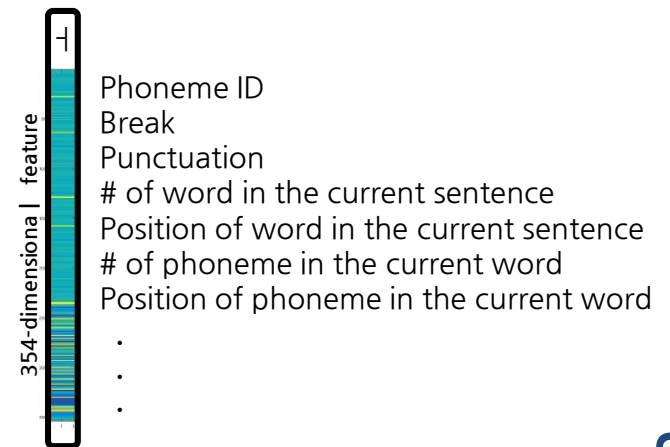
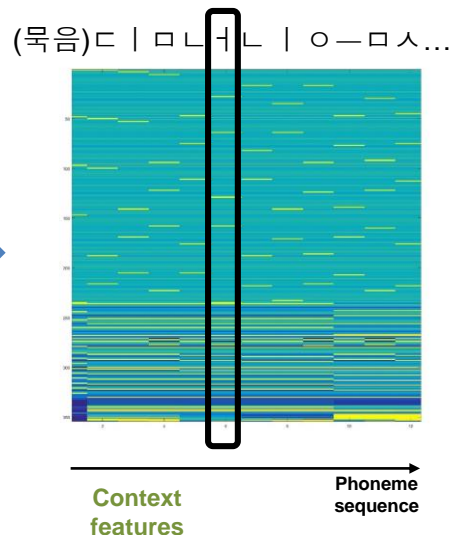
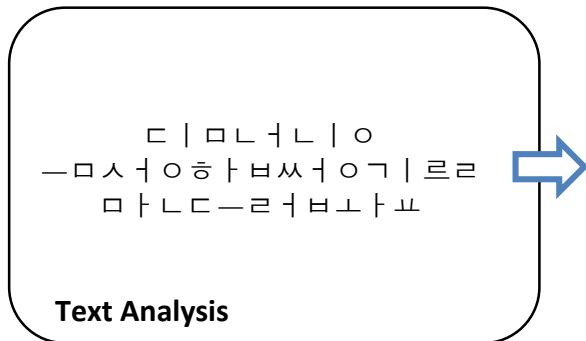
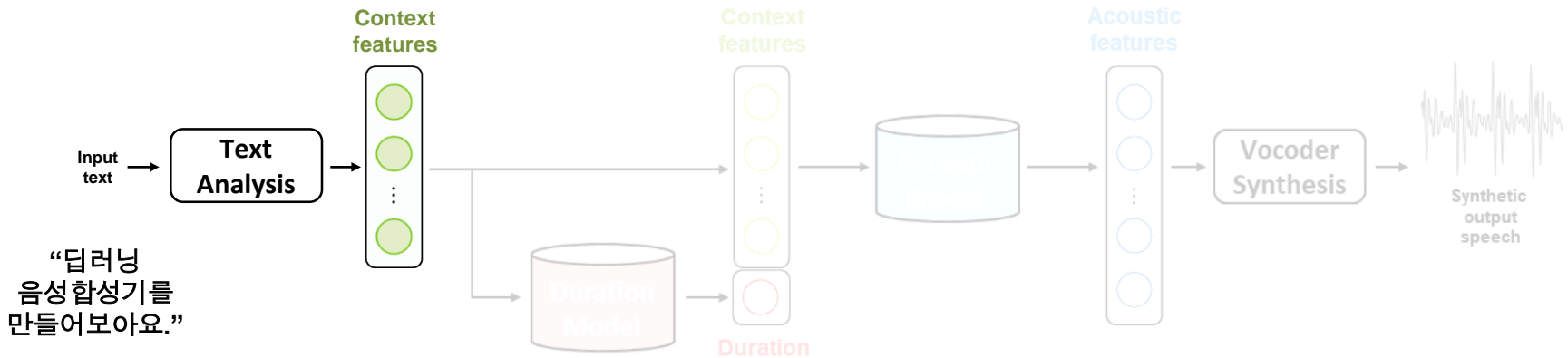


Context features

Phoneme sequence

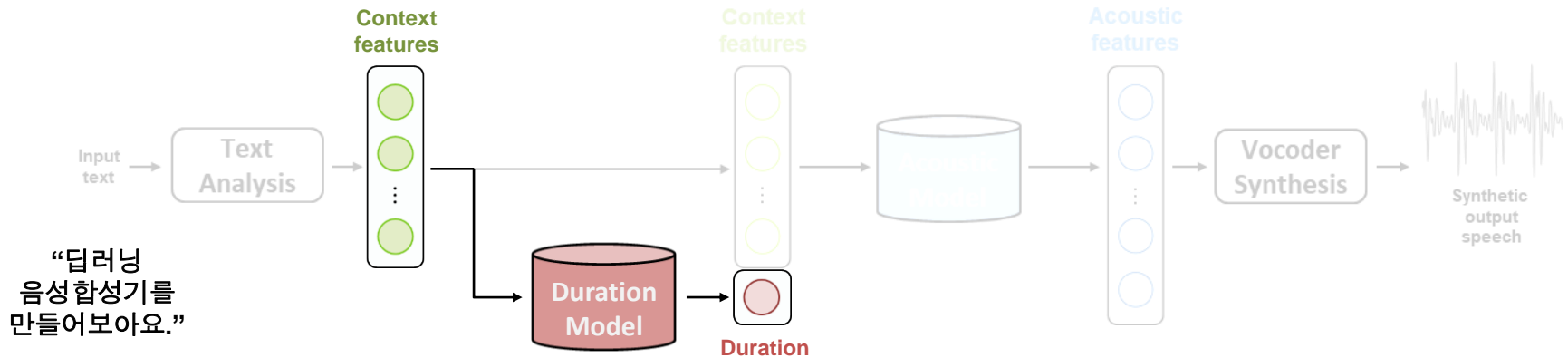
# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Text analyzer

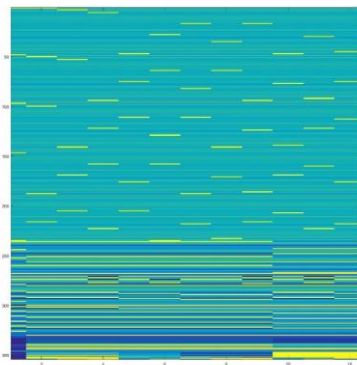


# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Duration model



(목음)ㄷ | 모ㄴ ㄱㄴ | ㅇ-ㅁㅅ...

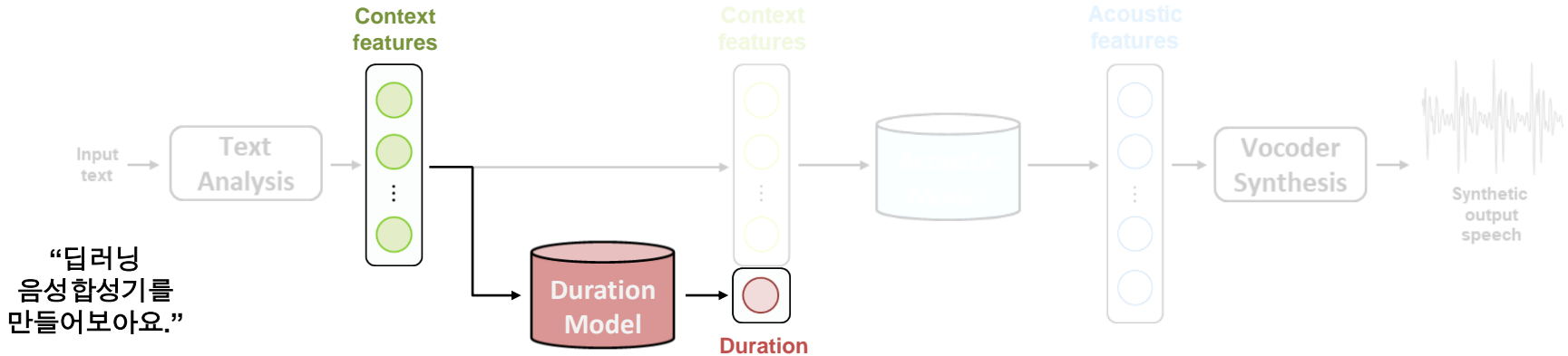


Context features

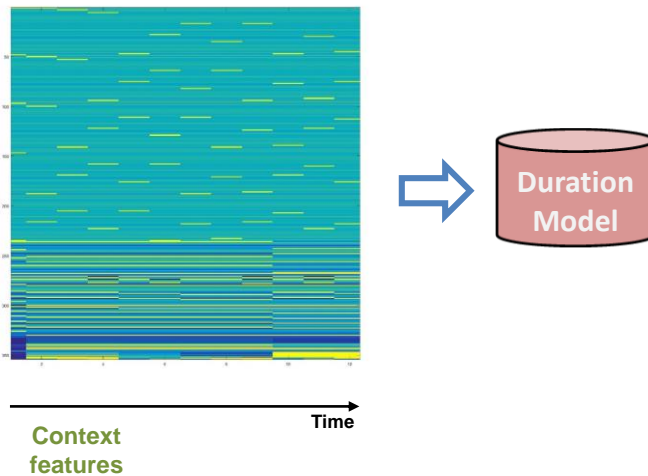
Phoneme sequence

# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Duration model

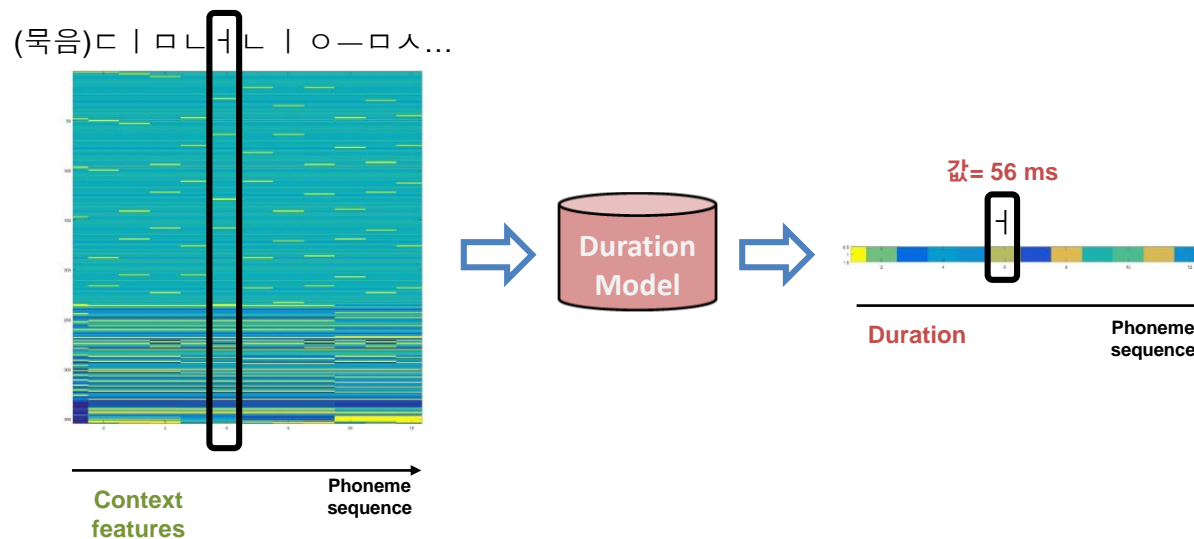
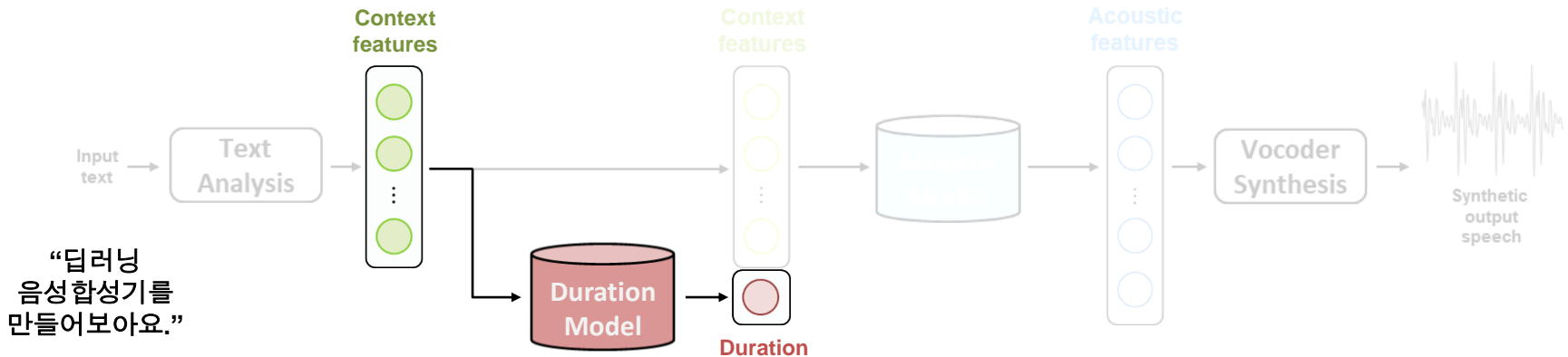


(목음)ㄷ | 모ㄴ ㄱㄴ | ㅇ-ㅁㅂ...



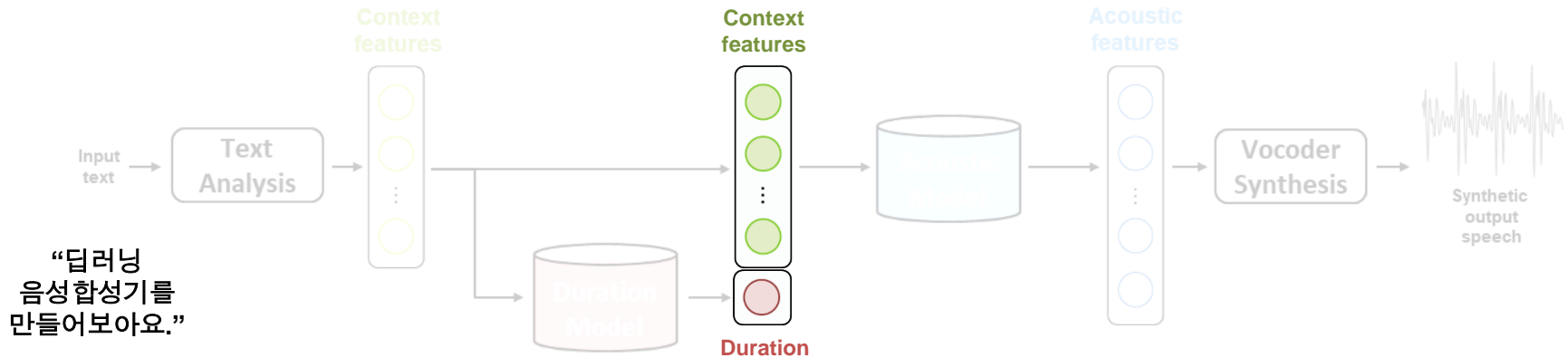
# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Duration model



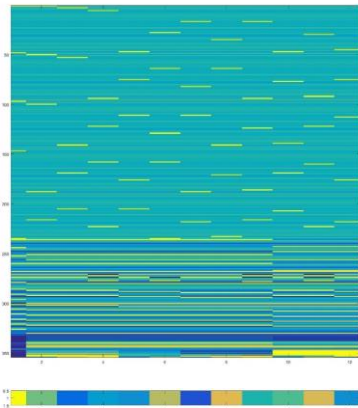
# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Linguistic upsampler



“딤러닝 음성합성기를 만들어보아요.”

(목음)ㄷ | ㄹㄴ | ㄷ | ㄴ | ㅇ - ㄹㅅ...

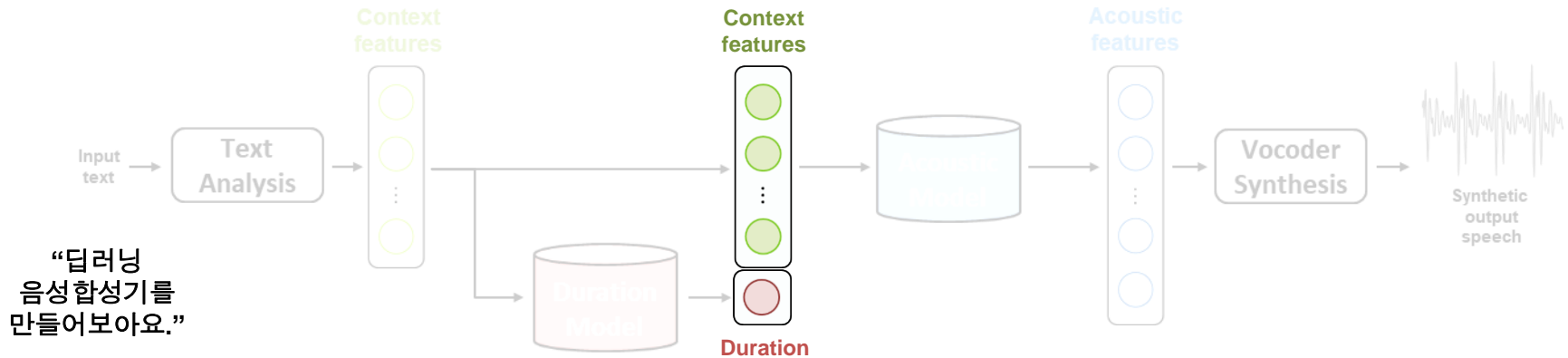


Context features + Duration Phoneme sequence

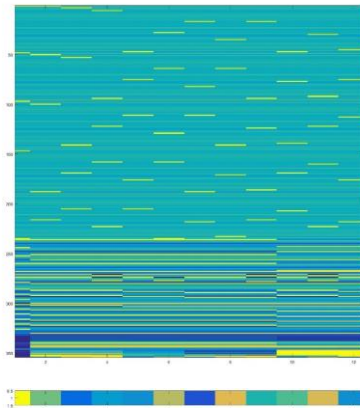


# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Linguistic upsampler



(목음)ㄷ | ㄹㄴ | ㄷㄴ | ㅇ-ㄹㅅ...

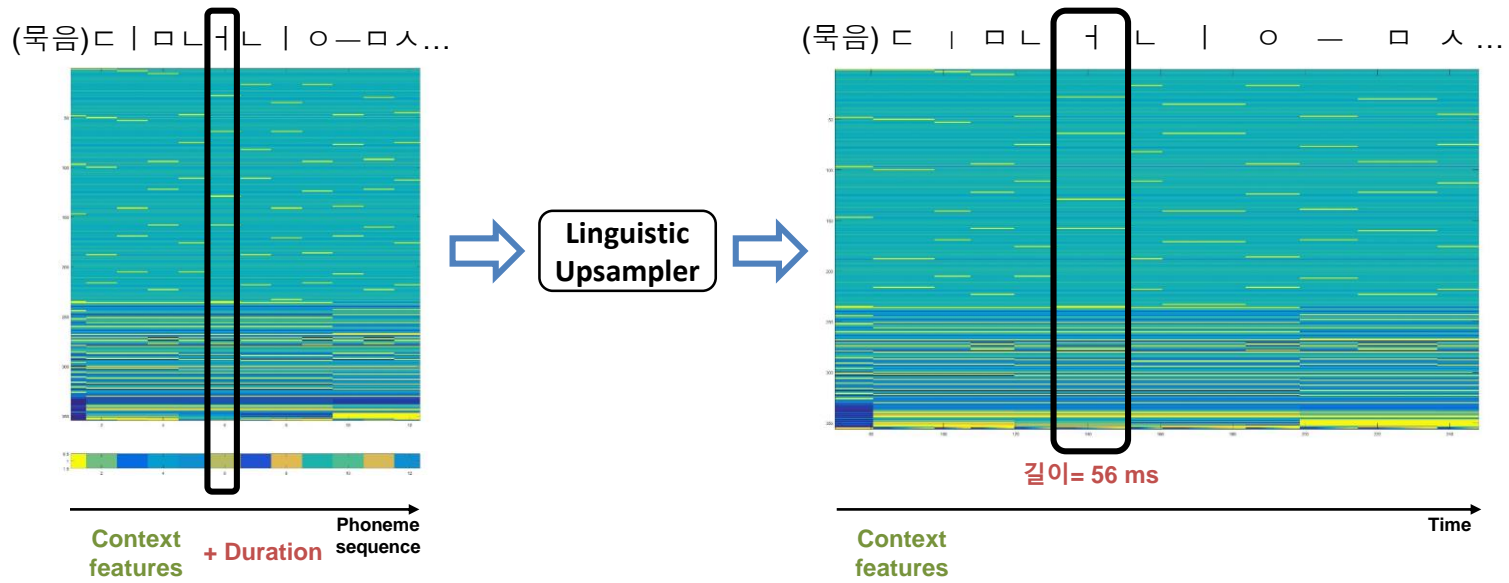
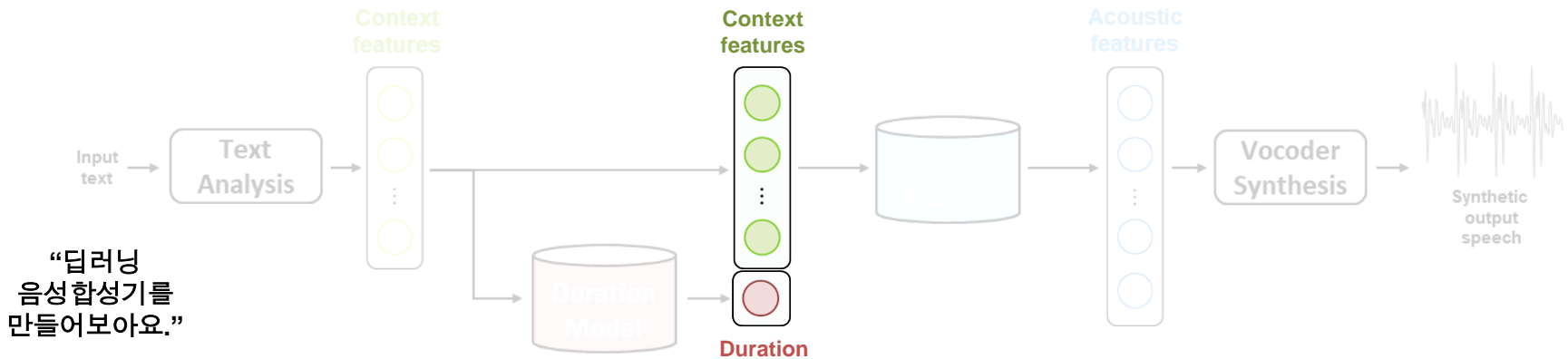


➡ Linguistic Upsampler

Context features + Duration  
Phoneme sequence

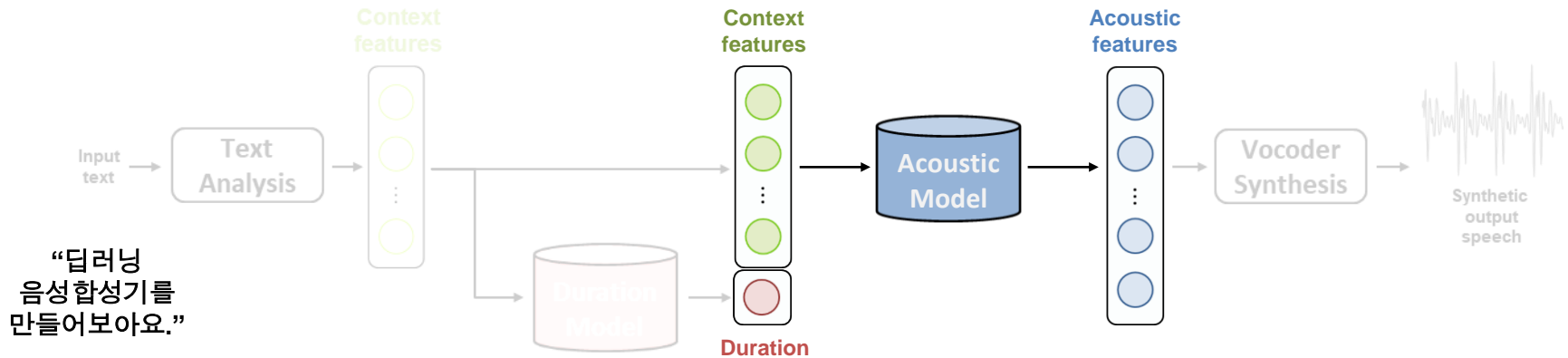
# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Linguistic upsampler



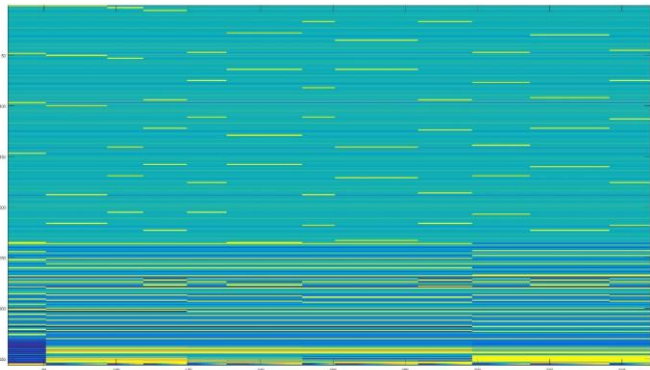
# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Acoustic model



“딥러닝 음성합성기를 만들어보아요.”

(목음) ㄷ | ㅁ ㄴ ㅅ ㄴ | ㅇ - ㅁ ㅅ ...

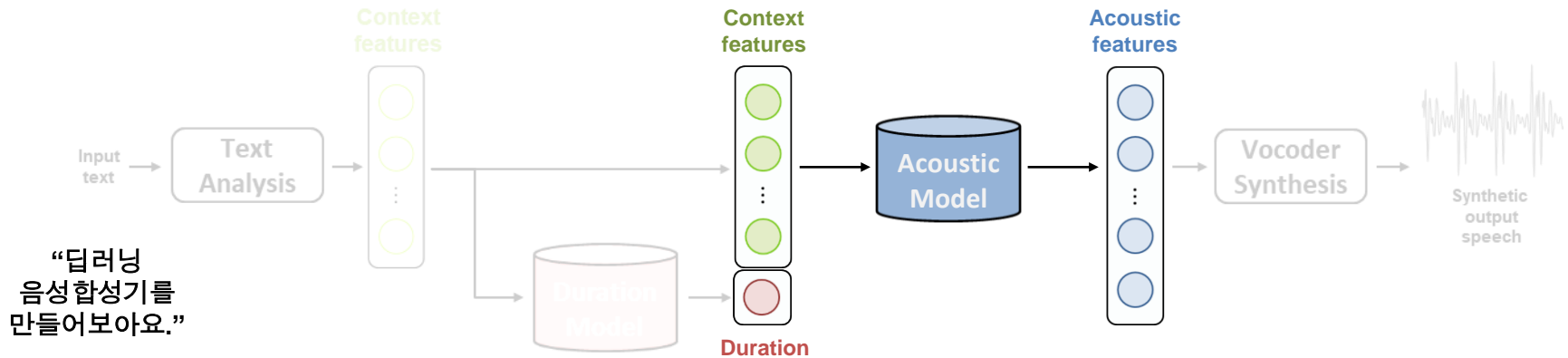


Context features

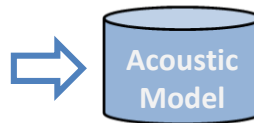
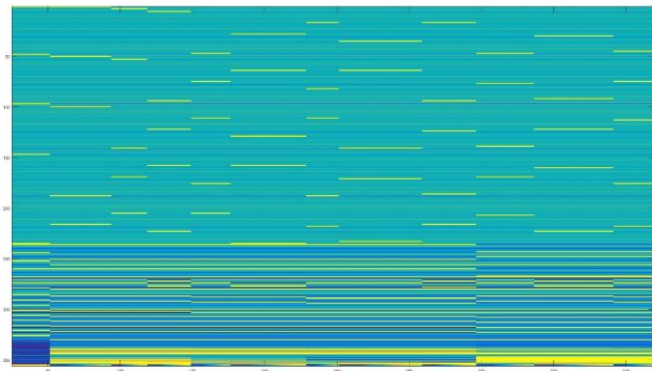
Time

# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Acoustic model



(목음) □ | □ L | L | ○ - □ s ...

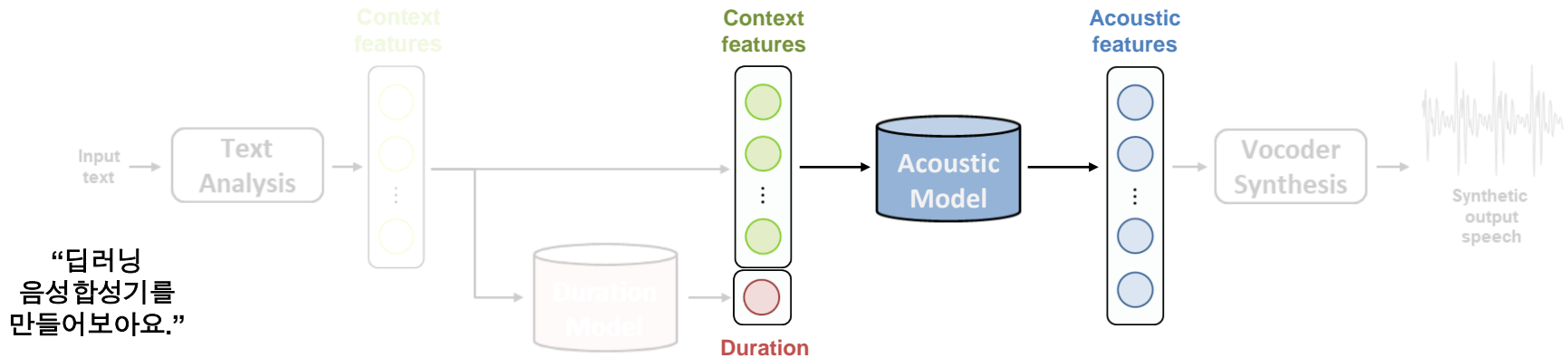


Context features

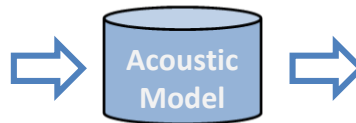
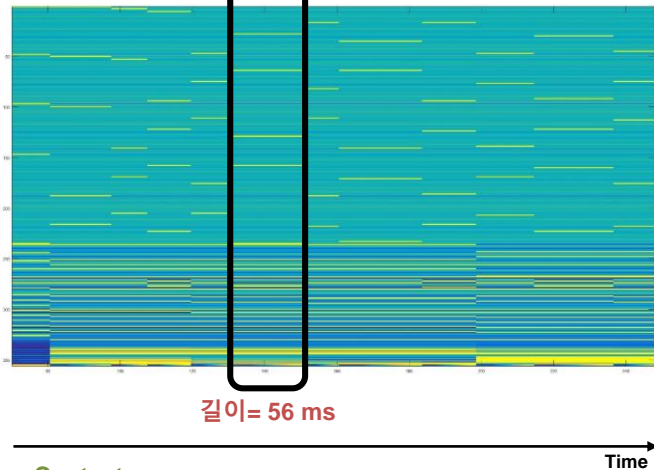
Time

# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

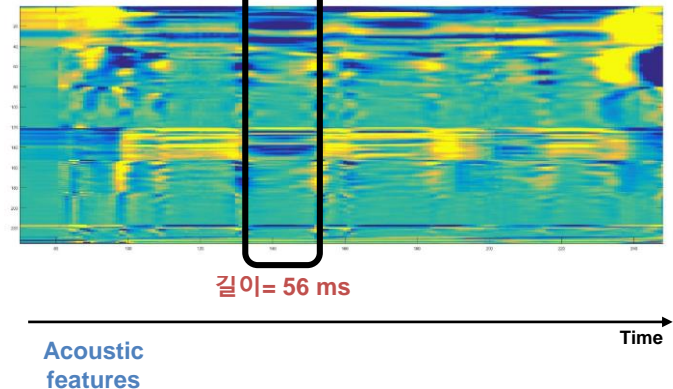
## Acoustic model



(목음) □ | □ □ □ | □ | □ □ □ □ □ □ □ ...

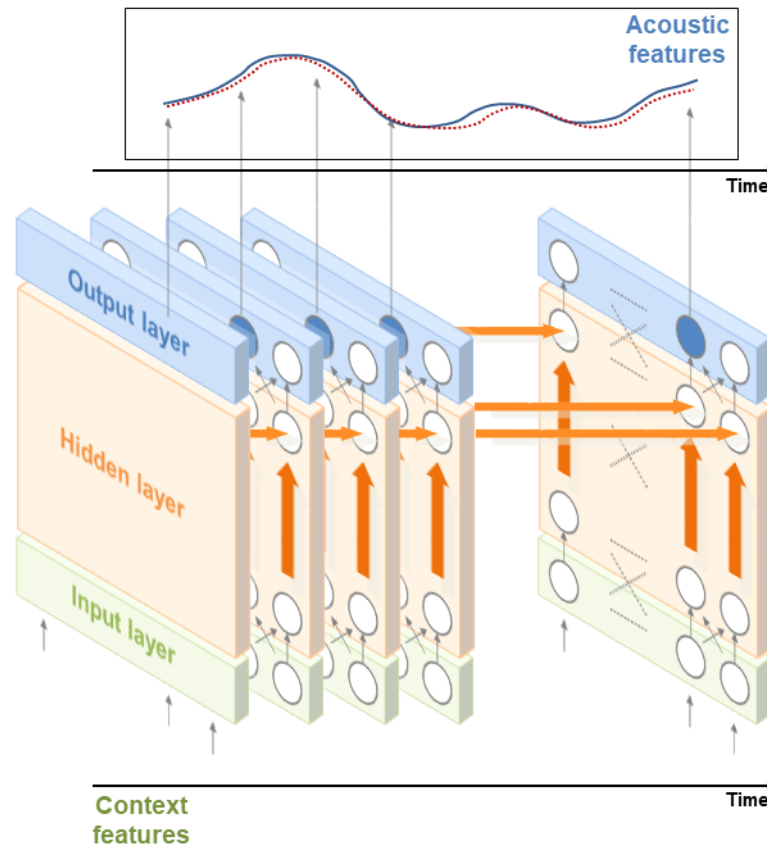


(목음) □ | □ □ □ | □ | □ □ □ □ □ □ □ ...



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

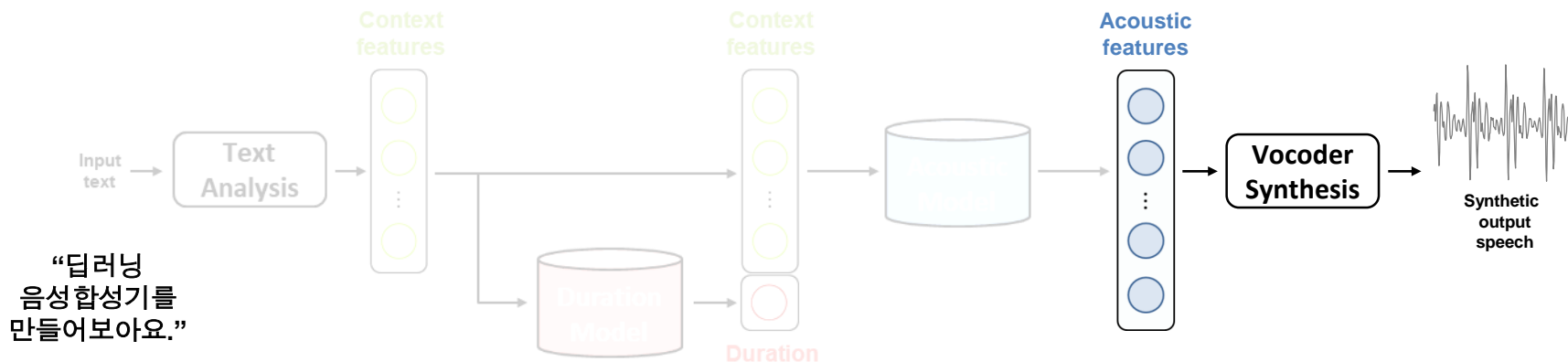
Multiple feed-forward and LSTM layers



$$\text{AcousticFeature} = F_{\text{DNN}}(\text{ContextFeature})$$

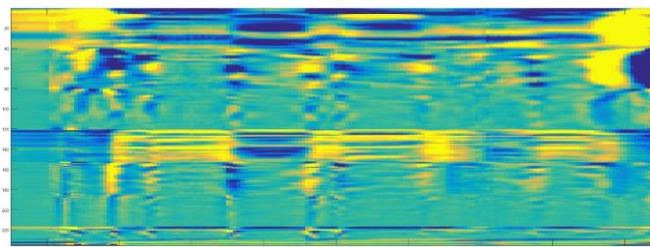
# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

## Vocoder synthesis



“딥러닝  
음성합성기를  
만들어보아요.”

(목음) □ | □ □ □ □ □ | ○ — □ □ ...



**Vocoder  
Synthesis**

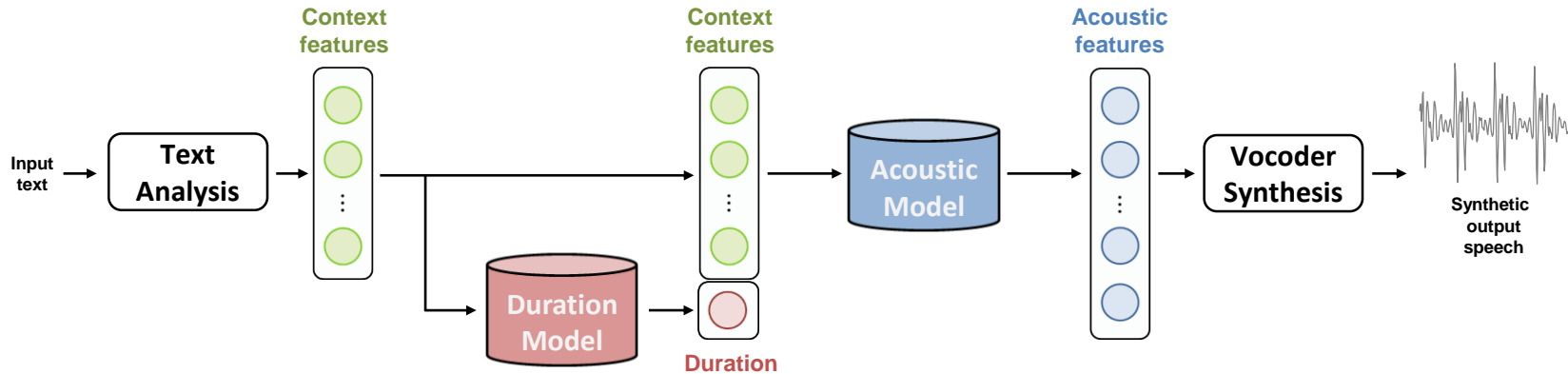


Acoustic  
features

Time

# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Generate speech parameters from input text and deep neural network

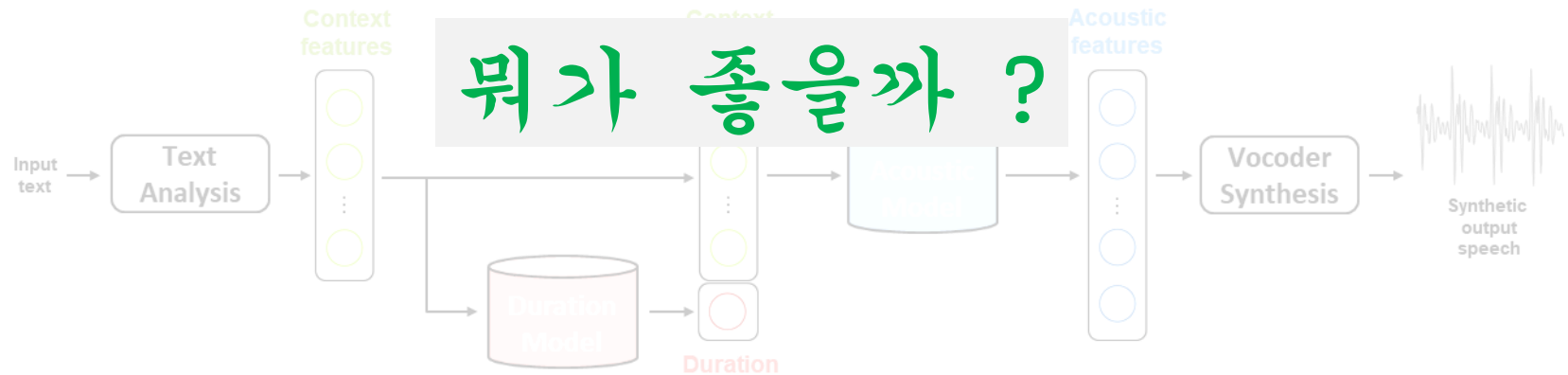


- NLP frontend : text analyzer
- Speech backend : vocoder synthesizer
- Deep learning model : duration & acoustic model



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Generate speech parameters from input text and deep neural network



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Generate speech parameters from input text and deep neural network



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Generate speech parameters from input text and deep neural network



뭐가 좋을까 ?

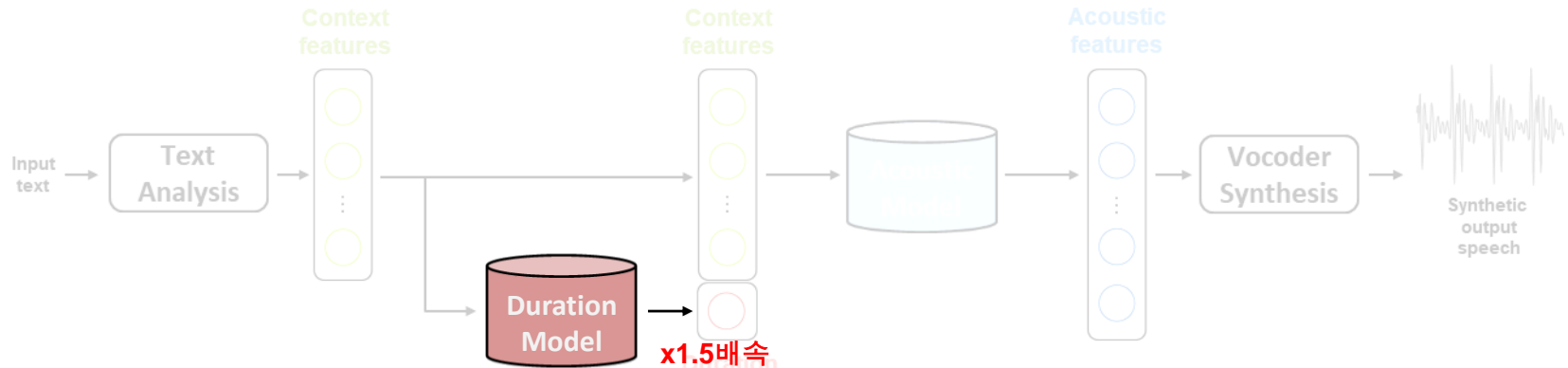
통제 모델을 사용하니까 ...

**Flexible to change  
voice characteristics**



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Applications : speech rate



간장 공장 공장장은  
강 공장장이고,  
된장 공장 공장장은  
공 공장장이다.

내가 그린 기린 그림은  
긴 기린 그림이고,  
네가 그린 기린 그림은  
안 긴 기린 그림이다.

좌로인정 우로인정  
앞구르기인정  
인정올리지말고 인정내려  
인정안해서 후회한다면  
후회 할 시간을  
후회하는 각이고요.

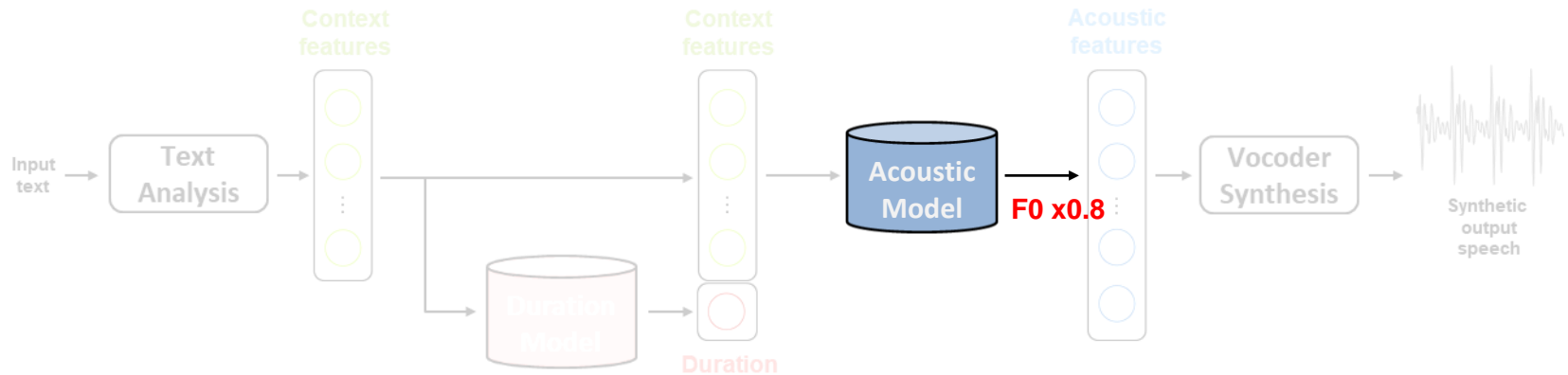
오졌따리 오저따  
콩콩따리콩콩따  
산기슭이  
인정하는 바이고요  
슬픔발이  
인정하는 바입니다.

x1.5배속



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Applications : speech tone



간장 공장 공장장은  
강 공장장이고,  
된장 공장 공장장은  
공 공장장이다.



내가 그린 기린 그림은  
긴 기린 그림이고,  
네가 그린 기린 그림은  
안 긴 기린 그림이다.



좌로인정 우로인정  
앞구르기인정  
인정올리지말고 인정내려  
인정안해서 후회한다면  
후회 할 시간을  
후회하는 각이고요.



오졌따리 오저따  
콩콩따리콩콩따  
산기슭이  
인정하는 바이고요  
슌곰발이  
인정하는 바입니다.

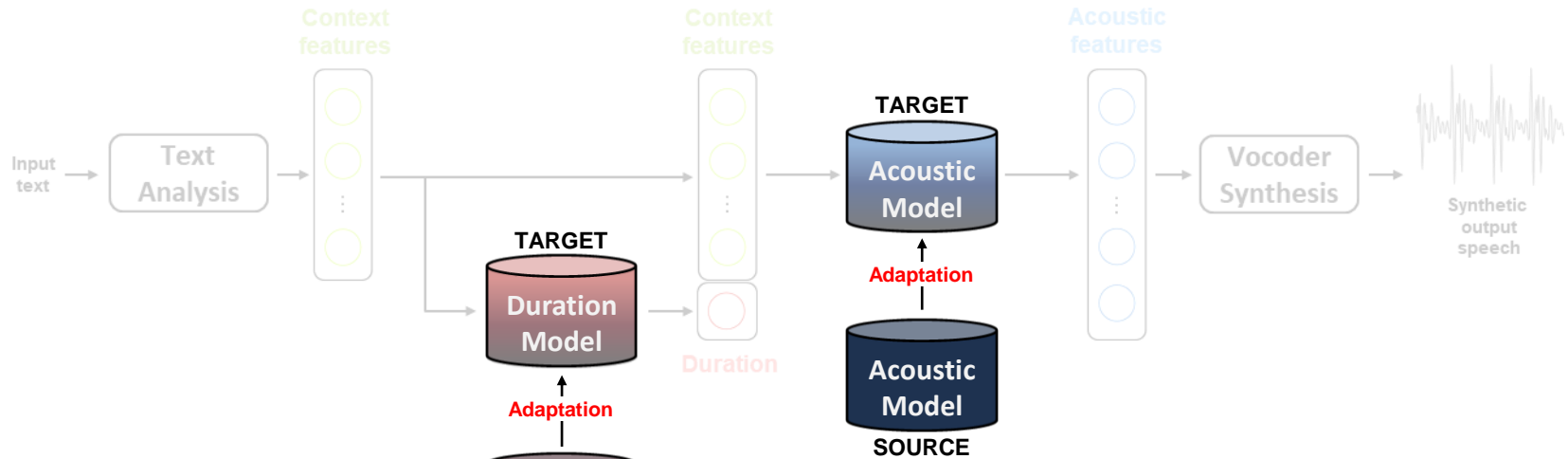


F0 x0.8



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Applications : speaker adaptation



간장 공장 공장장은  
강 공장장이고,  
된장 공장 공장장은  
공 공장장이다.

네가 그린 그림은  
안 긴 기린 그림이다.  
네가 그린 그림은  
안 긴 기린 그림이다.

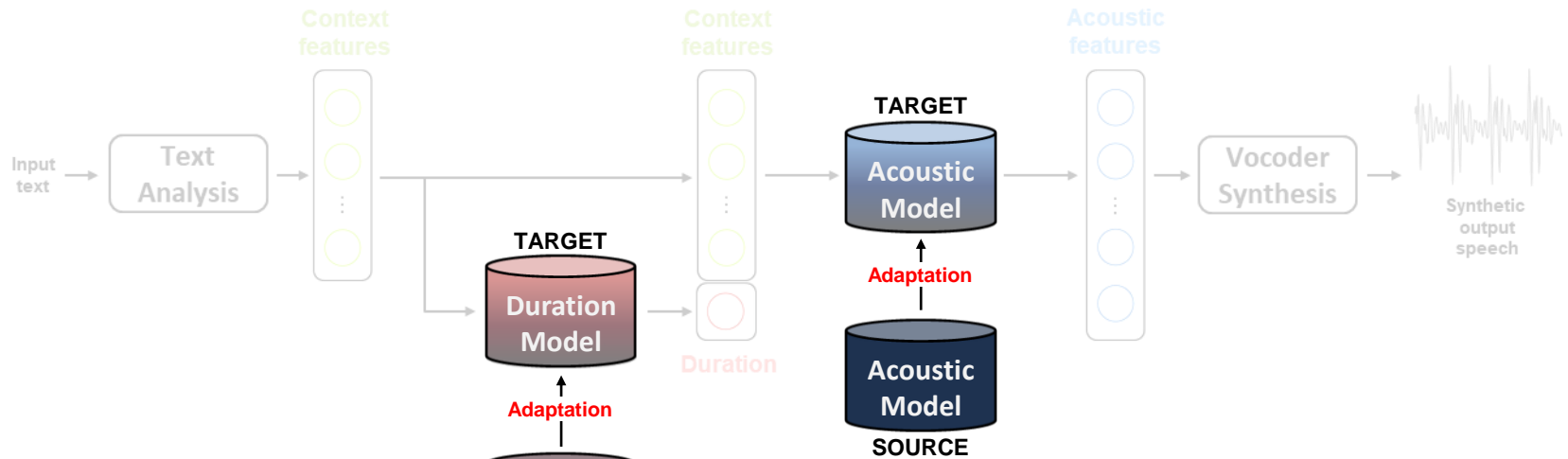
좌로인정 우로인정  
앞구르기인정  
인정올리지말고 인정내려  
인정안해서 후회한다면  
후회 할 시간을  
후회하는 각이고요.

오졌따리 오저따  
콩콩따리콩콩따  
산기슭이  
인정하는 바이고요  
슬픔발이  
인정하는 바입니다.



# DEEP-LEARNING SPEECH SYNTHESIS SYSTEM (DTS)

Applications : speaker adaptation



간장 공장 공장장은  
강 공장장이고,  
된장 공장 공장장은  
공 공장장이다.

네가 그린 그림은  
안 긴 기린 그림이다.  
네가 그린 기린 그림은  
안 긴 기린 그림이다.

좌로인정 우로인정      오졌따리 오져따  
후회하는 각이고요.      인정하는 바입니다.

**유인나 4시간**



x1.5배속



F0 x0.8



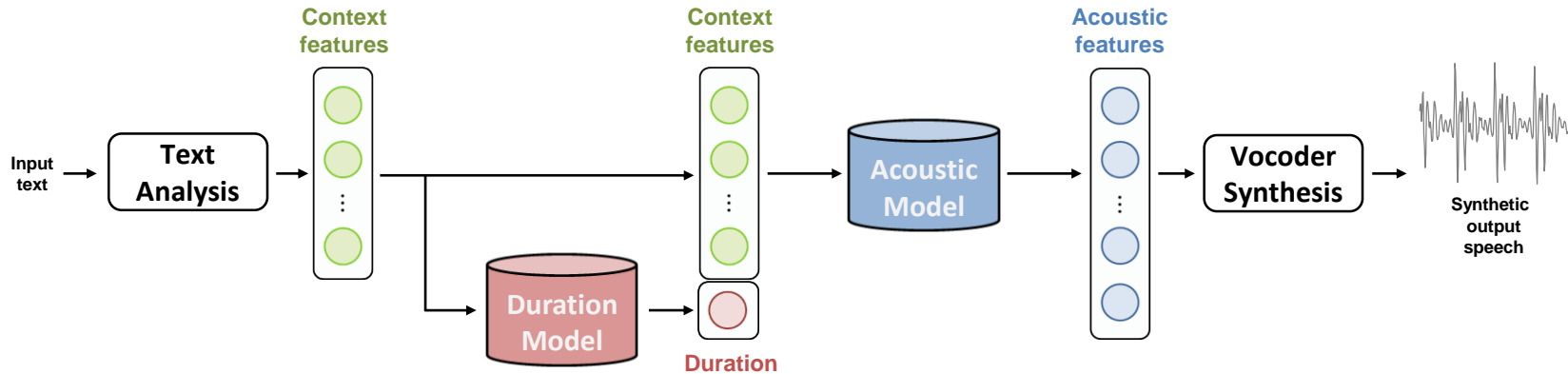
x1.5배속



F0 x0.8

# SUMMARY

Generate speech parameters from input text and deep neural network



- **NLP** frontend : text analyzer
- **Speech** backend : vocoder synthesizer
- **Deep learning** model : duration & acoustic model
  
- Advantages
  - Flexible to change voice characteristics



질문있어요 ?

CLOVA 채용부스로 오세요.

