

Audio & Speech Technology

[6] Speech Synthesis & Recognition

Phoneme

Phoneme = Smallest segment of speech sound

EXAMPLE

Vowels

Consonants

IPA	Example
i:	<u>see</u> , <u>heat</u>
a:	<u>arm</u> , <u>father</u>
u:	<u>blue</u> , <u>foot</u>
ɔ:	<u>call</u> , <u>horse</u>
ɜ:	<u>turn</u> , <u>learn</u>
ɪ	<u>ship</u> , <u>hit</u>
æ	<u>hat</u> , <u>black</u>
ʊ	<u>put</u> , <u>could</u>
ʌ	<u>cup</u> , <u>duck</u>

IPA	Example
b	<u>book</u>
d	<u>day</u> , <u>did</u>
f	<u>find</u> , <u>if</u>
g	<u>give</u> , <u>flag</u>
h	<u>how</u>
j	<u>yes</u> , <u>yellow</u>
k	<u>cat</u> , <u>back</u>

IPA	Example
l	<u>leg</u> , <u>little</u>
z	<u>zoo</u> , <u>lazy</u>
ʒ	<u>pleasure</u>
ʃ	<u>she</u> , <u>crash</u>
tʃ	<u>check</u>
n	<u>no</u> , <u>ten</u>
ŋ	<u>sing</u>

IPA = International Phonetic Alphabet

IPA = International Phonetic Alphabet

Thai Phonemes

SHORT VOWELS		LONG VOWELS		DIPHTHONGS	
IPA	ตัวอักษร	IPA	ตัวอักษร	IPA	ตัวอักษร
a	ะ, ั	a:	า, ำ	iaʔ, iəʔ	เียะ
e	เะ, เ็	e:	เ, เอ	ia, iə	เีย, เีย
ɛ	แะ, แ็	ɛ:	แ, แอ	uaʔ, uəʔ	ัวะ
i	อิ, ิ	i:	อี, ี	ua, uə	ัว, ำ
o	โะ, ็อ	o:	โ, เอ	waʔ, wəʔ	เือะ
ɔ	เาะ, ็อ	ɔ:	อ, ็อ	wa, wə	เือ, เือ
u	อุ, ุ	u:	ู, ู		(เสียงสระควบ)
w	อื, ุ	w:	ือ, ุ		
ʉ	เอะ	ʉ:	เอ, ุ		

(เสียงสระสั้น)

(เสียงสระยาว)

See more at http://en.wikipedia.org/wiki/International_Phonetic_Alphabet_chart_for_English_dialects

Thai Phonemes

CONSONANTS				TONES	
IPA	ตัวอักษร	IPA	ตัวอักษร	IPA	วรรณยุกต์
b	บ	p	ป	a	เสียงสามัญ
d	ฎ, ด	ph	ผ, พ, ภ	à	เสียงเอก
f	ฝ, ฟ	r	ร, ฬ	â	เสียงโท
h	ห, ฮ	s	ซ, ศ, ษ, ส	á	เสียงตรี
j	ญ, ย, อย, หย	t	ฏ, ต	ǎ	เสียงจัตวา
k	ก	th	ฐ, ท, ฒ, ถ, ฑ, ฒ		
kh (/x/)	ข, ฃ, ค, ฅ, ฆ	tc	จ		
l	ล, ฬ, หล	tcʰ	ฉ, ช, ฌ		
m	ม, หม	w	ว, หว		
n	ณ, น, หน	ʔ	อ, ะ		
ŋ	ง, หง				

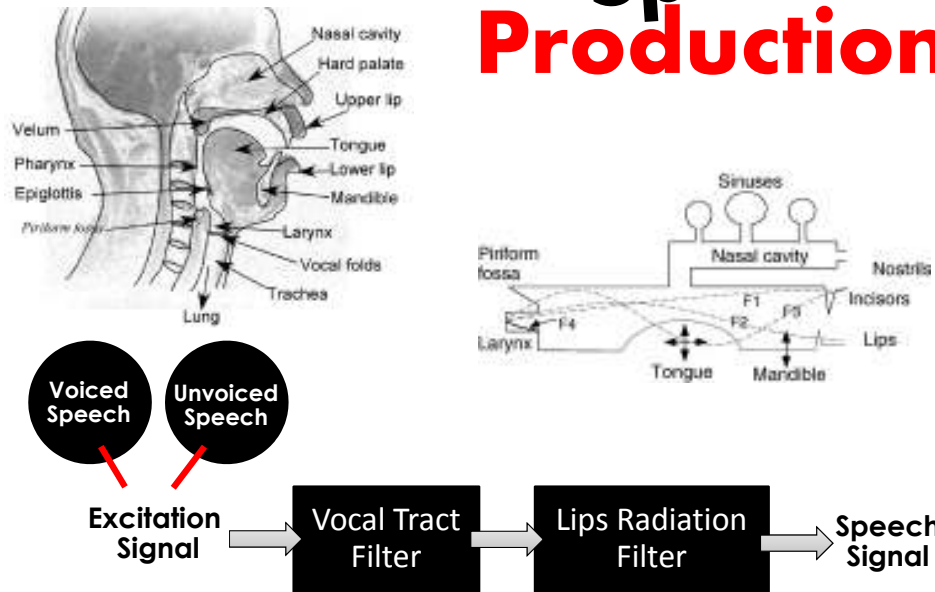
Ex

เจ็บมาก



ญั:bmak

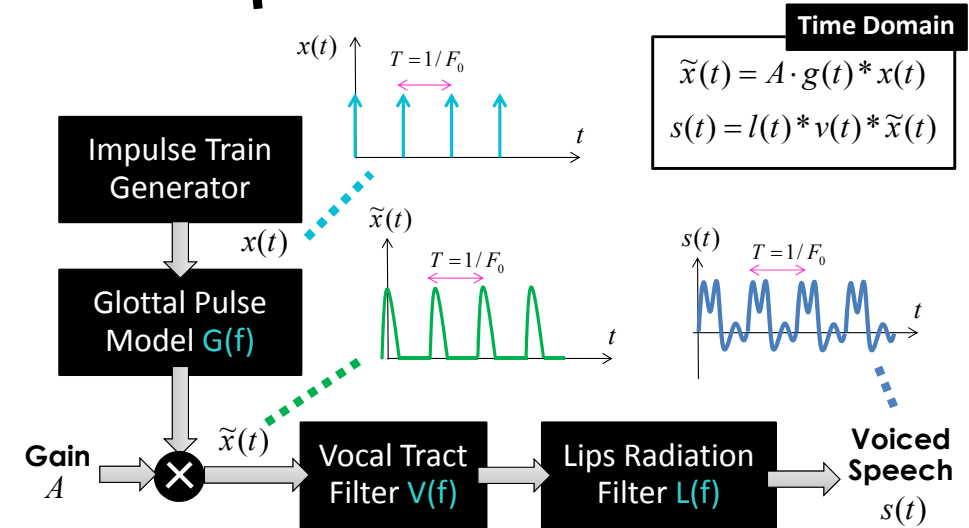
Speech Production



<http://ars.els-cdn.com/content/image/1-s2.0-S016763930700177X-gr1.jpg>

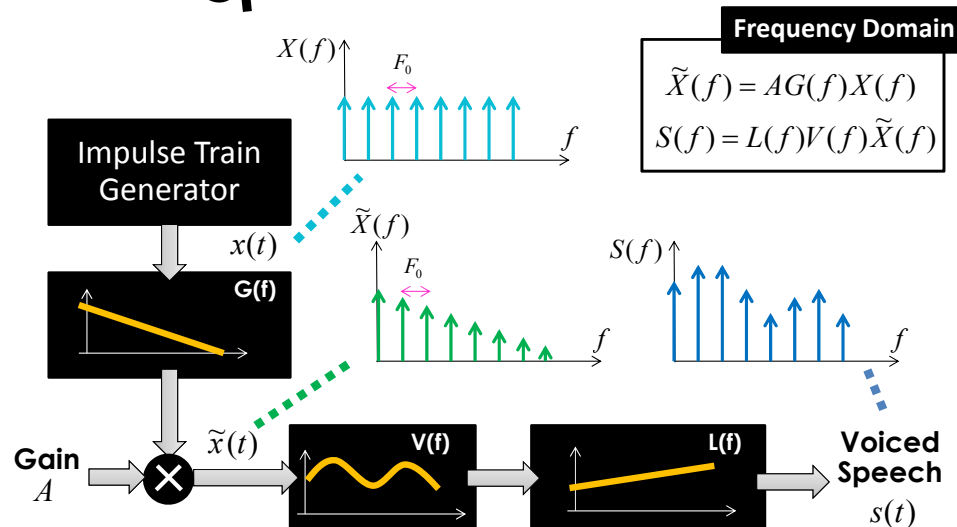
Voiced Speech

- Vocal cords vibrate
- Excitation is periodic signal
- Characterized by Low Freq.



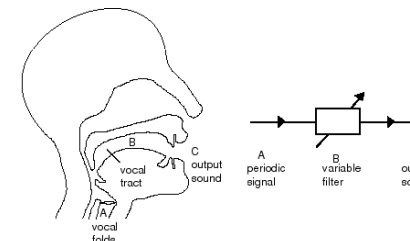
Voiced Speech

- Vocal cords vibrate
- Excitation is periodic signal
- Characterized by Low Freq.



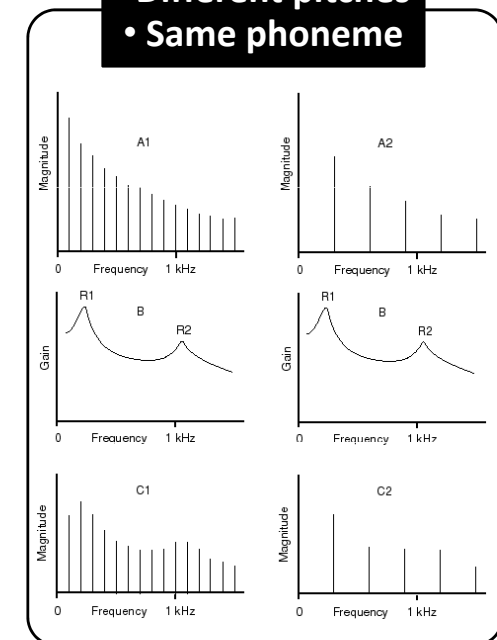
Voiced Speech

- Different pitches
- Same phoneme



Fundamental Frequency

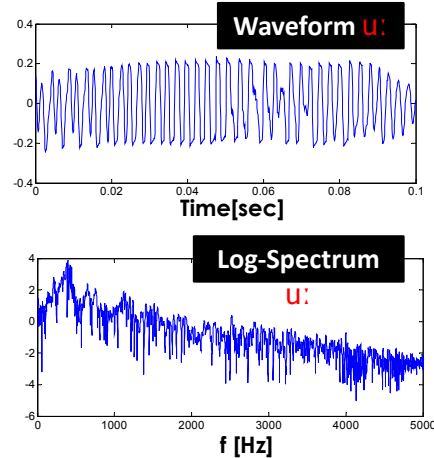
- Male : 50-200 Hz
- Female : 150-300 Hz
- Child : 200-400 Hz



Voiced Speech

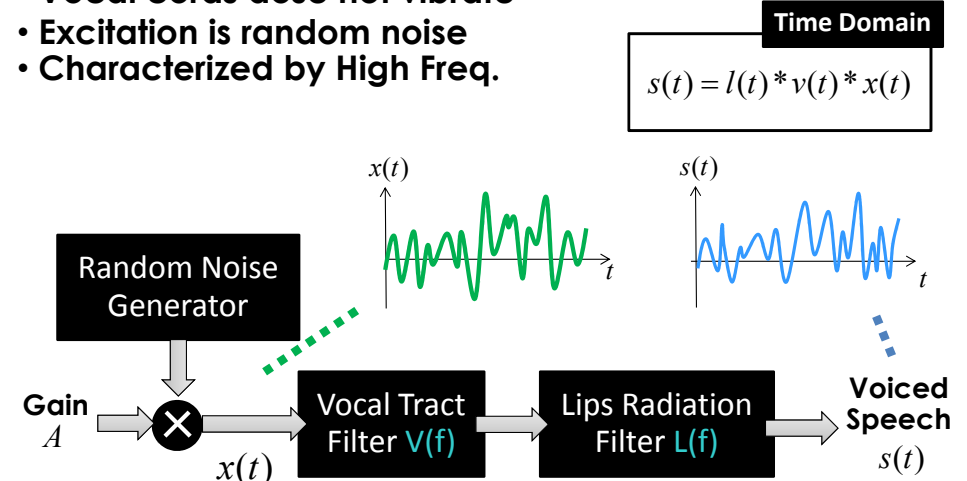
- All vowel in English = Voiced speech.
- Voiced consonants

IPA	Example	IPA	Example
b	<u>book</u>	l	<u>leg</u> , <u>little</u>
d	<u>day</u> , <u>did</u>	m	<u>man</u>
g	<u>give</u> , <u>flag</u>	n	<u>no</u> , <u>ten</u>
v	<u>vanilla</u>	ŋ	<u>sing</u>
z	<u>zoo</u> , <u>lazy</u>	r	<u>red</u> , <u>try</u>
ð	<u>then</u>	w	<u>window</u>
ʒ	<u>pleasure</u>	j	<u>yes</u> , <u>yellow</u>
dʒ	<u>jump</u> , <u>gin</u>		



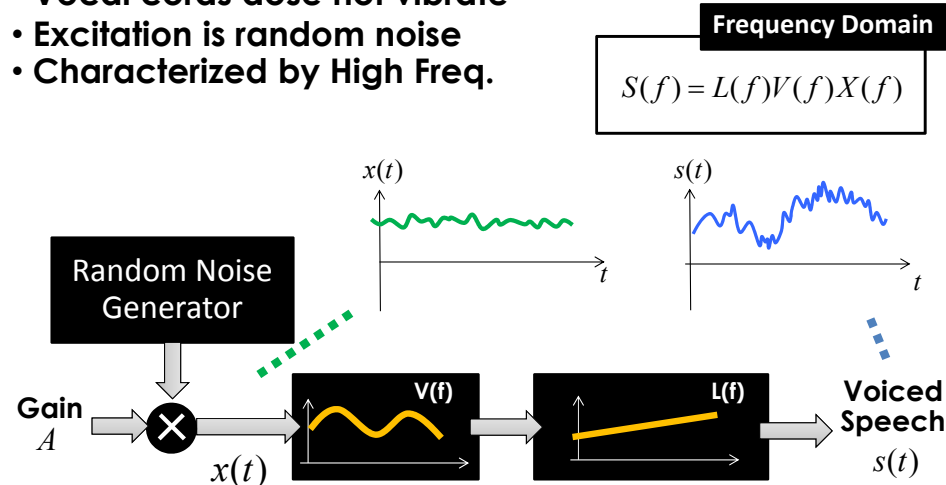
Unvoiced Speech

- Vocal cords dose not vibrate
- Excitation is random noise
- Characterized by High Freq.



Unvoiced Speech

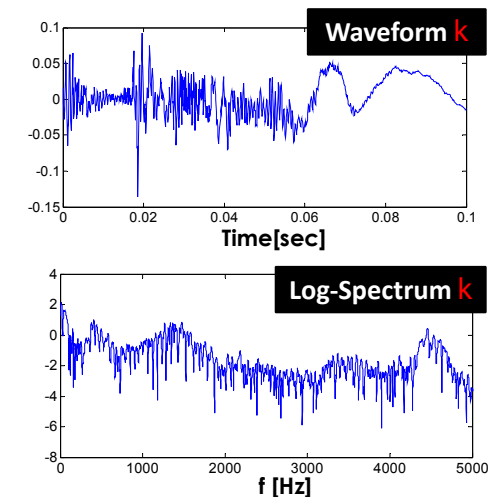
- Vocal cords dose not vibrate
- Excitation is random noise
- Characterized by High Freq.



Unvoiced Speech

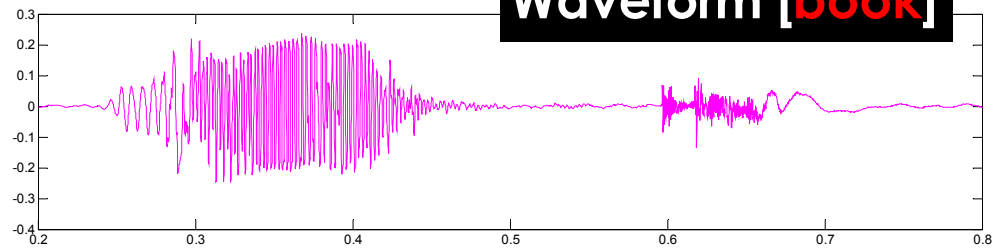
- Unvoiced consonants

IPA	Example
p	<u>pet</u>
f	<u>find</u> , <u>if</u>
θ	<u>thirty</u> , <u>both</u>
t	<u>ten</u>
s	<u>sir</u>
ʃ	<u>she</u> , <u>crash</u>
tʃ	<u>check</u>
k	<u>king</u> , <u>cat</u>

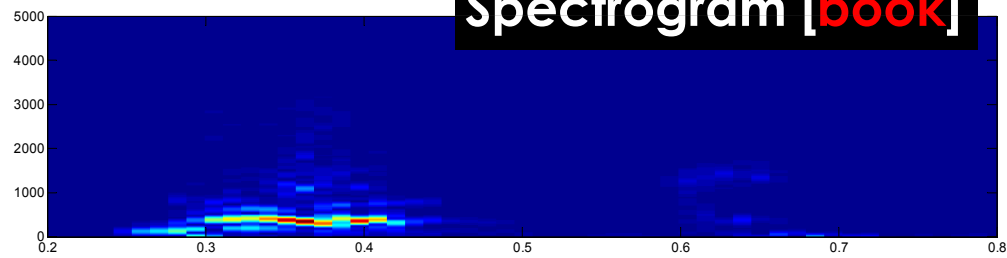


Voiced vs Unvoiced

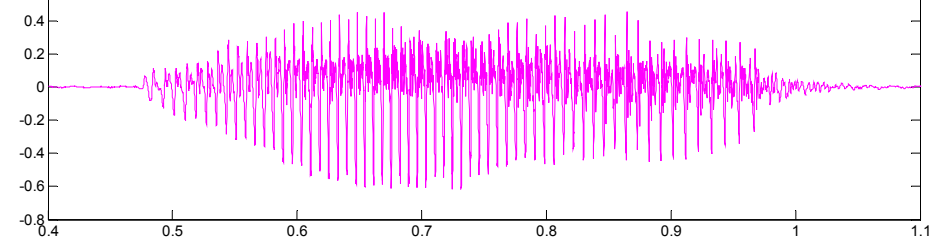
Waveform [book]



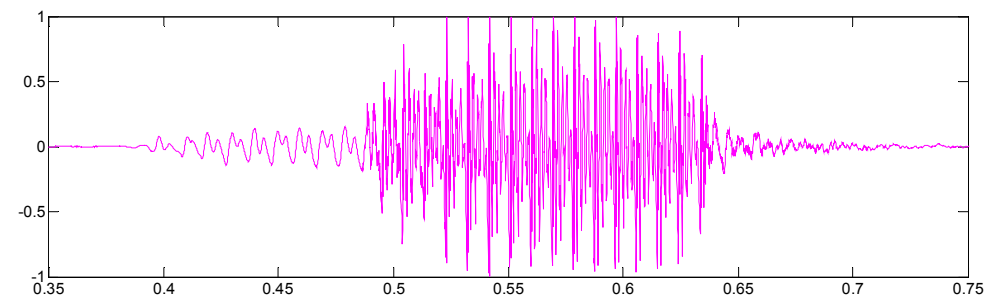
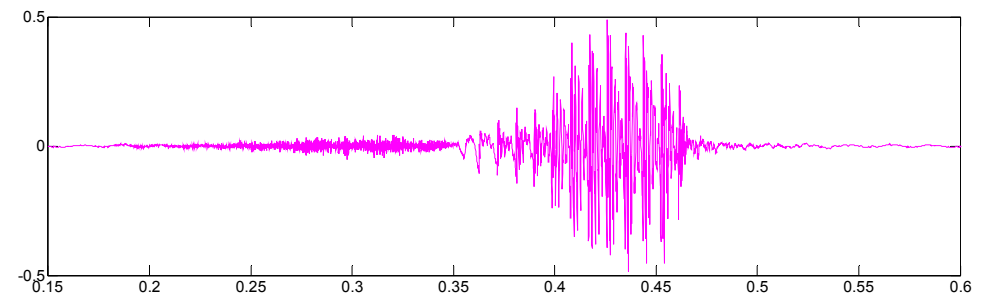
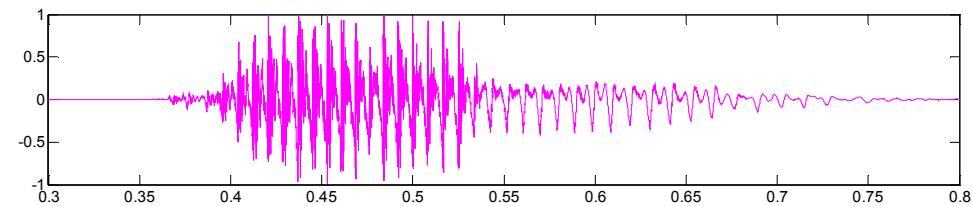
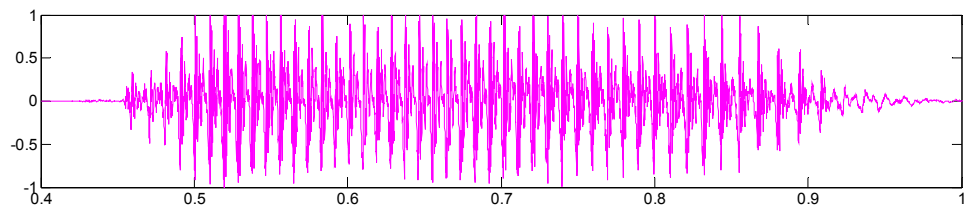
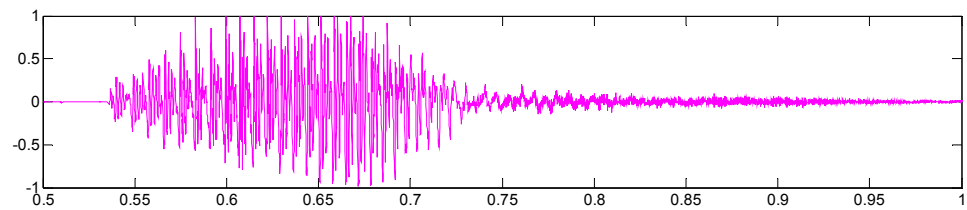
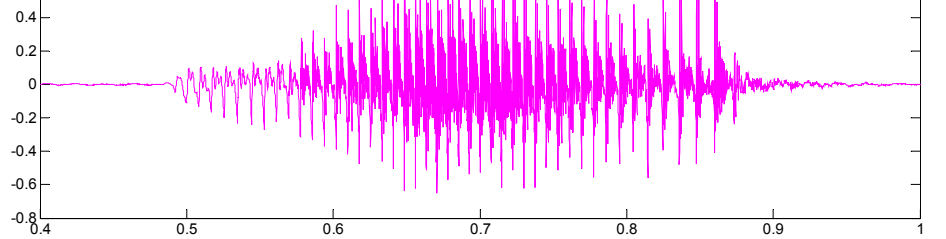
Spectrogram [book]



Waveform [เจ็บ]

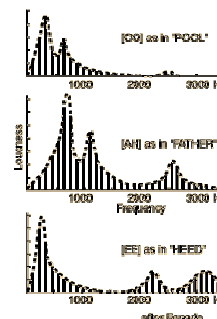


Waveform [งาบ]



Formant

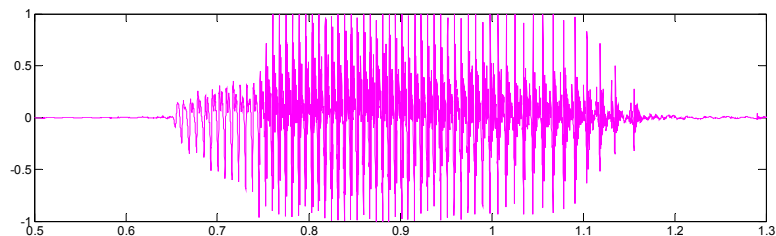
- **Vocal tract** = Resonance Tube
- **Formant** = Resonant Frequency of Vocal Tract
- Vocal tract has a fixed characteristic in the order of 10 ms
- 3-4 formants present below 4kHz of speech



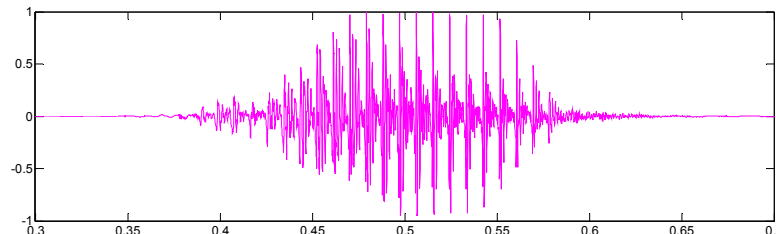
Phonetic Symbol	Example Word	F_1 (Hz)	F_2 (Hz)	F_3 (Hz)
/ow/	bought	570	840	2410
/oo/	boot	300	870	2240
/u/	foot	440	1020	2240
/a/	hot	730	1090	2440
/uh/	but	520	1190	2390
/er/	bird	490	1350	1690
/ae/	bat	660	1720	2410
/e/	bet	530	1840	2480
/i/	bit	390	1990	2550
/iy/	beet	270	2290	3010

<http://hyperphysics.phy-astr.gsu.edu/hbase/music/imgmus/vow5.gif>
http://cnx.org/content/m15459/latest/sub_formants-voweltable.png

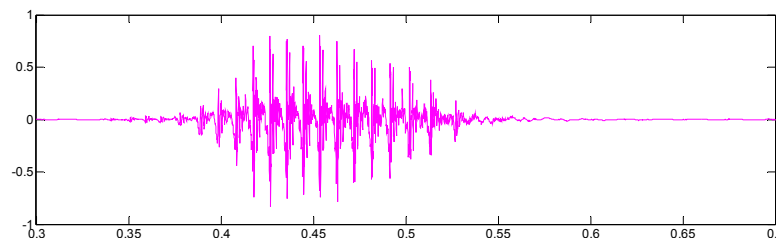
A



B

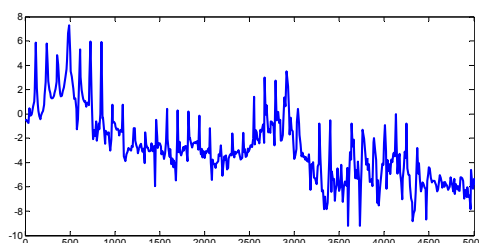
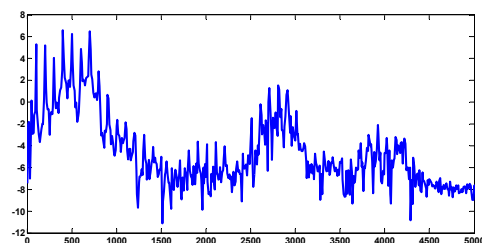


C

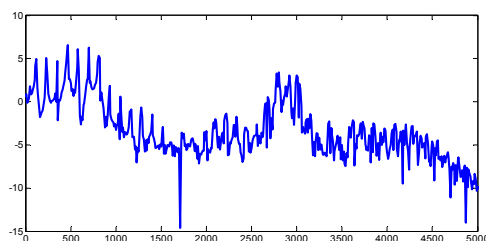


Formant

โอดาวว์

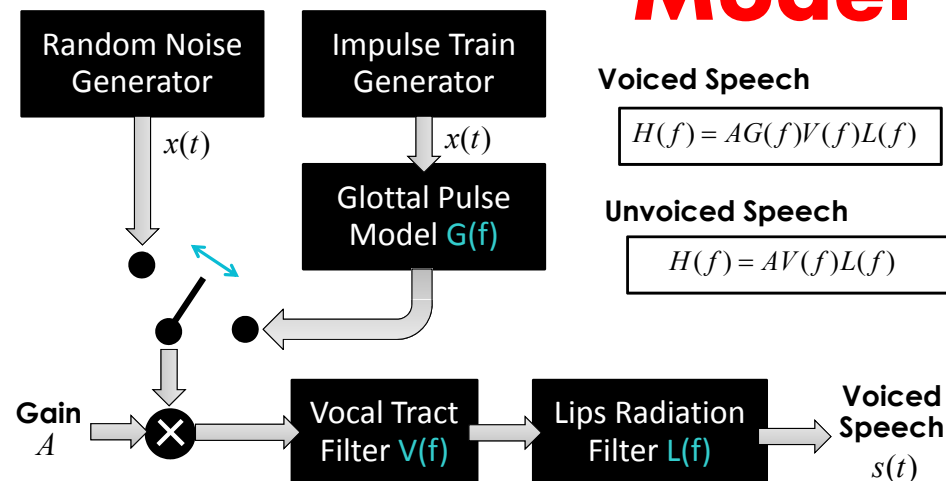


Formants = ?



Speech Production Model

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



Frequency Domain

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

Speech Production Model

Time Domain

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

$K = \text{Volumn Control}$
 $a_i = \text{Filter Coefficients}$

Depend on phonemes and speakers
 Can be estimated by using LPC

$$P \in \{10, 12, 14\}$$

$s[n] = \text{Speech Signal}$

$x[n] = \text{Excitation}$

Impulse Train for voiced speech [Varying Pitch]

Random noise for unvoiced Speech

Linear Predictive Coding

Current sample

Previous samples

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

LPC Coefficients

Prediction Error

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

Linear Predictive Coding

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

Speech Waveform
[Phoneme]

Training
LPC Model

LPC Coefficients
 $\{a_i\}$

Speech Production Model

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

Linear Predictive Coding

Speech Waveform
[Unknown]

$s[n]$

LPC Model

Predicted Speech Waveform
 $\hat{s}[n]$

Prediction Error
 $\varepsilon[n]$



LPC Model

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

Speech Production Model

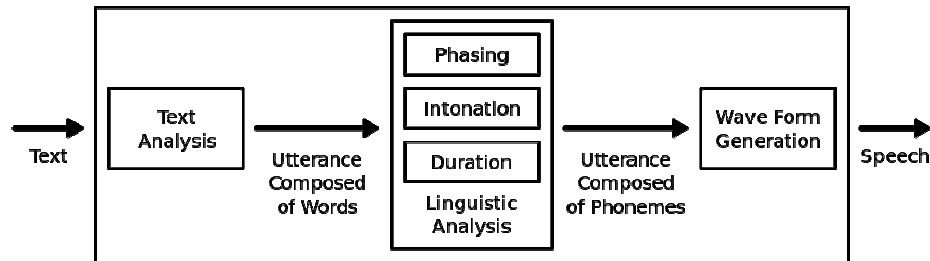
$$s[n] = \sum_{k=1}^P a_k s[n-k] + K \cdot x[n]$$

Excitation can be estimated from prediction error

- Voiced/Unvoiced Classes
- Pitches

$$\varepsilon[n] = K \cdot x[n]$$

Speech Synthesis

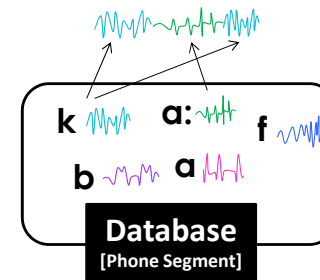


TTS : Text-to-Speech

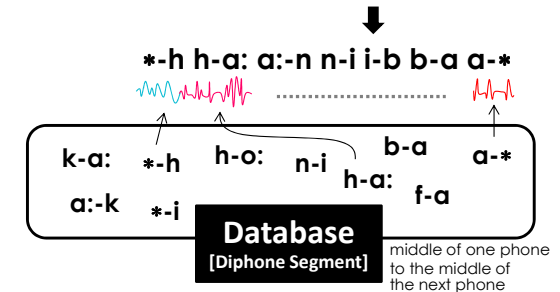
http://en.wikipedia.org/wiki/Speech_synthesis

Concatenative Synthesis

กาก → k a: k



ฮานิบะ → h a: n i b a



- Concatenate segments of **recorded speech**
- Produce natural synthesized speech
- Large databases of segmented recorded speech (phones, diphones, words, phrases, sentences)
- Unit selection synthesis, Diphone synthesis, Domain-specific synthesis

Time & Pitch Manipulation

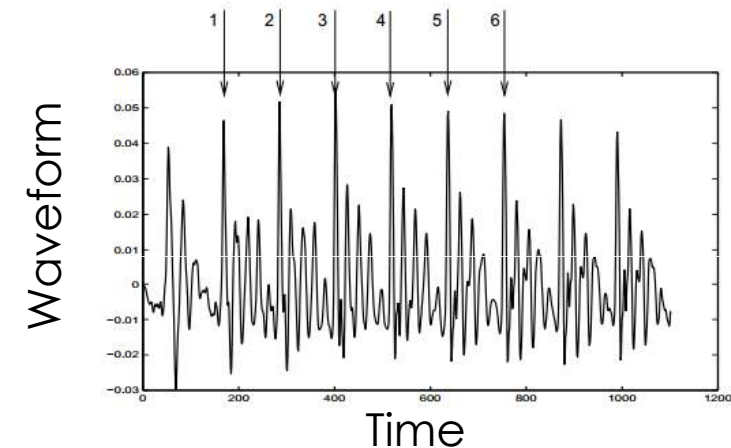
- **Time Stretching** = Change duration of sound (Shorten/Lengthen) without affecting its pitches
- **Pitch shifting/Scaling** = Change pitch of sound (Lower/Higher) without affecting its duration

Techniques

- Resampling
- Phase Vocoder
- PSOLA (Pitch Synchronous Overlap and Add)

PSOLA Pitch Synchronous Overlap and Add

Epoch (Pitch Mark) = Single instant in each pitch period that serves as an "anchor"



Formant Synthesis

Single Formant Synthesis Filter

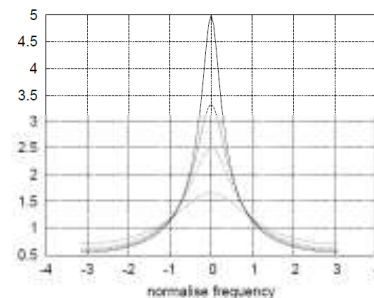
$$y[n] = s[n] + 2e^{-\pi B/F_S} \cos(2\pi \frac{F_R}{F_S}) y[n-1] - e^{-2\pi B/F_S} y[n-2]$$

F_R = Resonance Frequency (Formant)

B = Bandwidth of Resonance Filter

F_S = Sampling Frequency

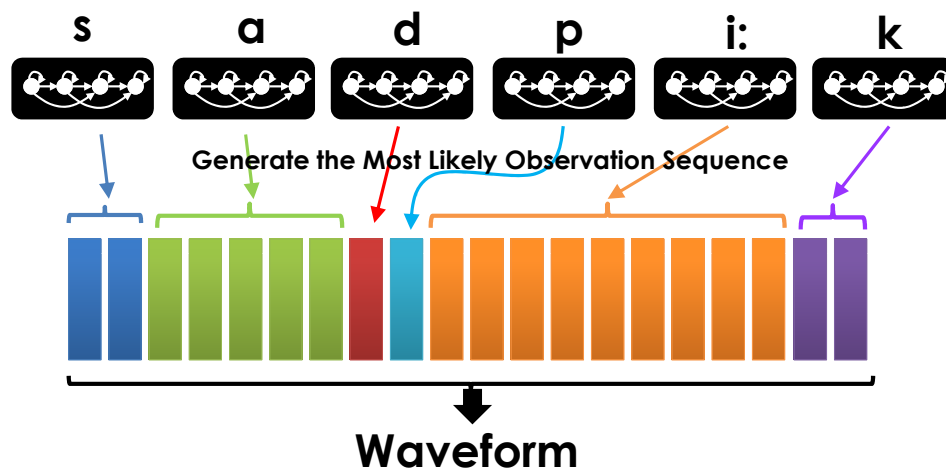
	/a/	/e/	/i/	/o/	/u/	BW
F1	750	469	281	468	312	90
F2	1187	2031	2281	781	1219	110
F3	2595	2687	3187	2656	2469	170
F4	3781	3375	3781	3281	3406	250
F5	4200	4200	4200	4200	4200	300



Paul Taylor, Text-to-Speech Synthesis
Alain de Cheveigné, FORMANT BANDWIDTH AFFECTS THE IDENTIFICATION OF COMPETING VOWELS

HMM-based Synthesis

สัทวีก → s a d p i: k
Duration 2 5 1 1 9 2 Frames



HMM-based Synthesis

We know **phone sequence** and **duration** for each phone. The phone sequence tells us **which models to use in which order**, but not which states to use, or which observations to generate. A duration tells us **how many observations that should be generated**, but again not which states to generate from.



Generate the **most likely sequence** of observations from the sequence of models

Each state will generate its mean observation

Speech Recognition

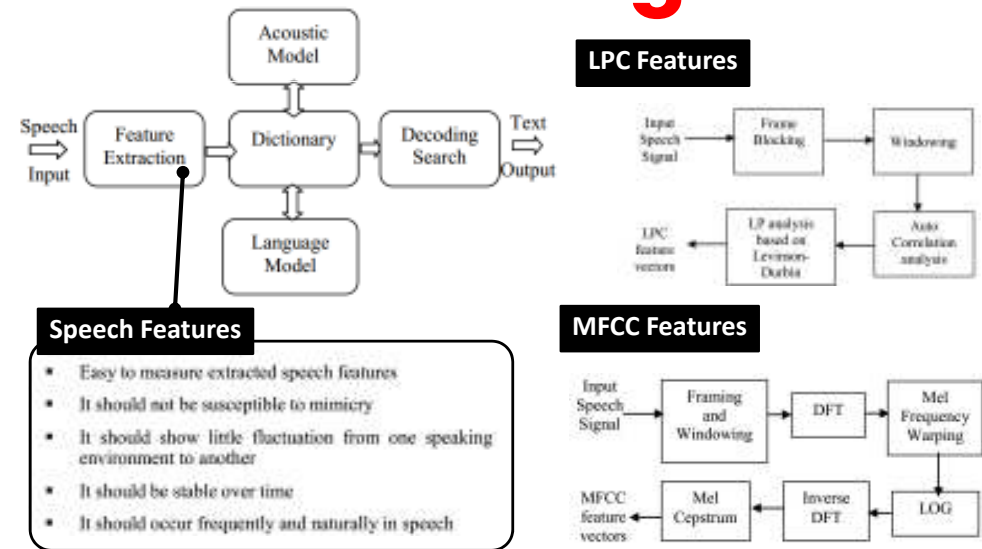
Speaker Model	Speaker Dependence [used by a single speaker] Speaker Independence [used by any speaker]
Vocabulary Size	Small Vocabulary [10 words] Medium Vocabulary [100 words] Large Vocabulary [1000 words] Very Large Vocabulary [10000 words] Out-of-Vocabulary [10000 words]
Speech Utterance	Isolated Words [single word] Connected Words / Discontinuous Speech [full sentence separated by silence] Continuous Speech [Naturally spoken sentences] Spontaneous Speech [Include mispronunciations, false-starts, and nonwords]

Growth of ASR

Year	Progress of ASR System
1952	Digit Recognizer
1976	1000 word connected recognizer with constrained grammar
1980	1000 word LSM recognizer (separate words w/o grammar)
1988	Phonetic typewriter
1993	Read texts (WSJ news)
1998	Broadcast news, telephone conversations
1998	Speech retrieval from broadcast news
2002	Rich transcription of meetings, Very Large Vocabulary, Limited Tasks, Controlled Environment
2004	Finnish online dictation, almost unlimited vocabulary based on morphemes
2006	Machine translation of broadcast speech
2008	Very Large Vocabulary, Limited Tasks, Arbitrary Environment
2009	Quick adaptation of synthesized voice by speech recognition (in a project where TKK participates in)
2011	Unlimited Vocabulary, Unlimited Tasks, Many Languages, Multilingual Systems for Multimodal Speech Enabled Devices
Future Direction	Real time recognition with 100% accuracy, all words that are intelligibly spoken by any person, independent of vocabulary size, noise, speaker characteristics or accent.

Vimala.C. & Dr.V.Radha, A Review on Speech Recognition Challenges and Approaches

Speech Recognition



Vimala.C. & Dr.V.Radha, A Review on Speech Recognition Challenges and Approaches

Speech Recognition

Approaches

- Hidden Markov Models (HMM)
 - Monophone Models
 - Triphone Models
- Neural Networks (NN)
- Dynamic Time Warping (DTW)

Vimala.C. & Dr.V.Radha, A Review on Speech Recognition Challenges and Approaches

HMM Speech Recognition

$$\hat{W} = \arg \max_w P(W | X) = \arg \max_w \frac{P(X | W)P(W)}{P(X)}$$

X = Observation Sequence ที่มี

W = Symbol Sequence ที่เป็นไปได้
Sequence of phoneme (word, phrase, sentences)

\hat{W} = Symbol Sequence ที่ Recognized ได้
สำหรับ Observation Sequence X

$P(X|W)$ = Prob. ของการเกิด X
จากค่า/วลี/ประโยค W

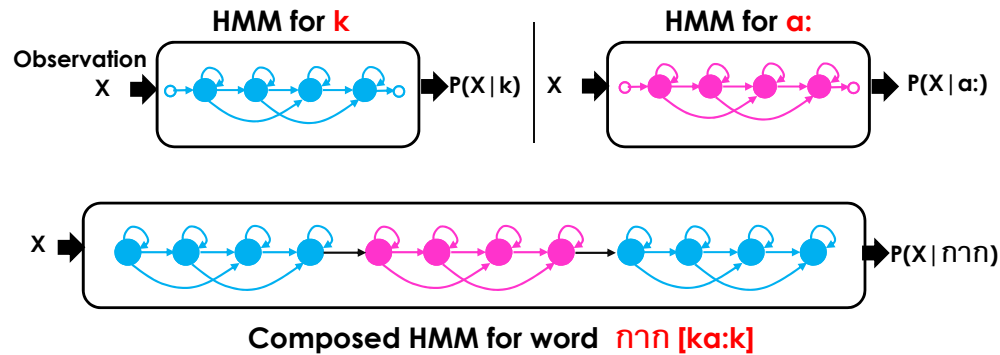
หาได้จาก Acoustic Model

$P(W)$ = Prob. ของการเกิด
ค่า/วลี/ประโยค W

หาได้จาก Language Model

Acoustic Model

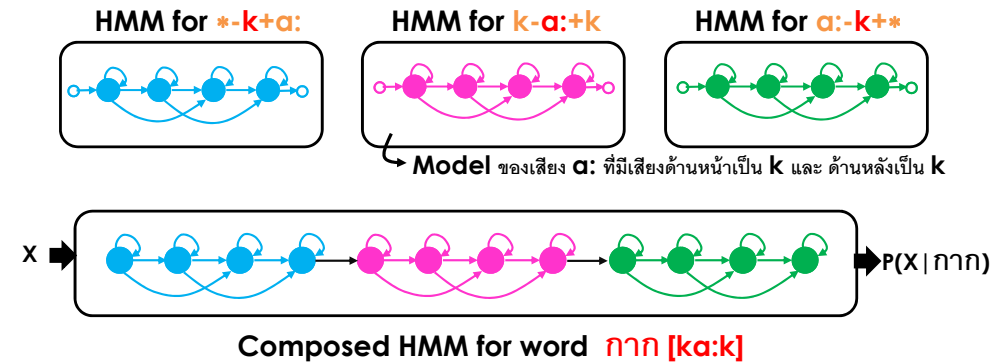
- Statistical model that describe the probability of acoustic observation sequence



Monophone Model
[Context Independent]

Acoustic Model

- Statistical model that describe the probability of acoustic observation sequence



Triphone Model
[Context Dependent]

Acoustic Model

Monophone Model [Context Independent]

แม่ม่ง = m, ϵ , η , o, η

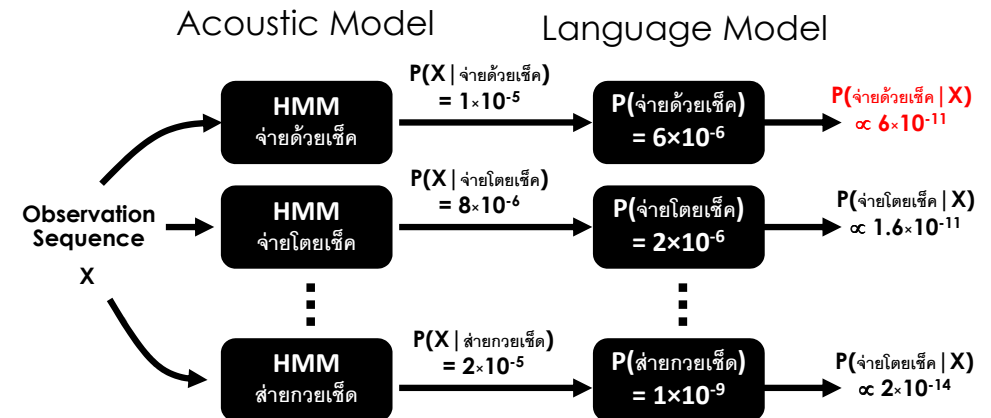
Triphone Model [Context Dependent]

Word Internal แม่ม่ง = $*-m+\epsilon$, $m-\epsilon+*$, $*-\eta+o$, $\eta-o+\eta$, $o-\eta+*$

Cross-Word แม่ม่ง = $*-m+\epsilon$, $m-\epsilon+\eta$, $\epsilon-\eta+o$, $\eta-o+\eta$, $o-\eta+*$

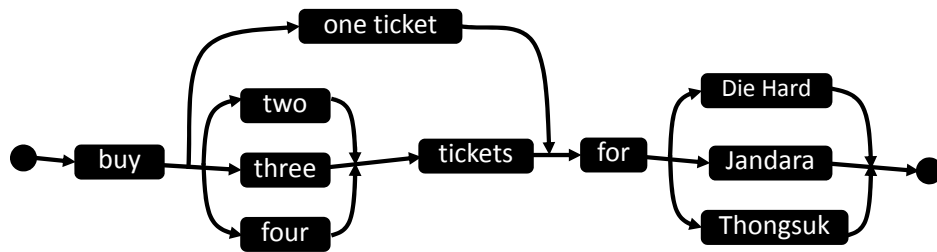
- Statistical Model that describe the probability of the symbol sequence (word/ phrase/ sentence)

Language Model



Language Model

Finite State Networks



Assume that words are Independent

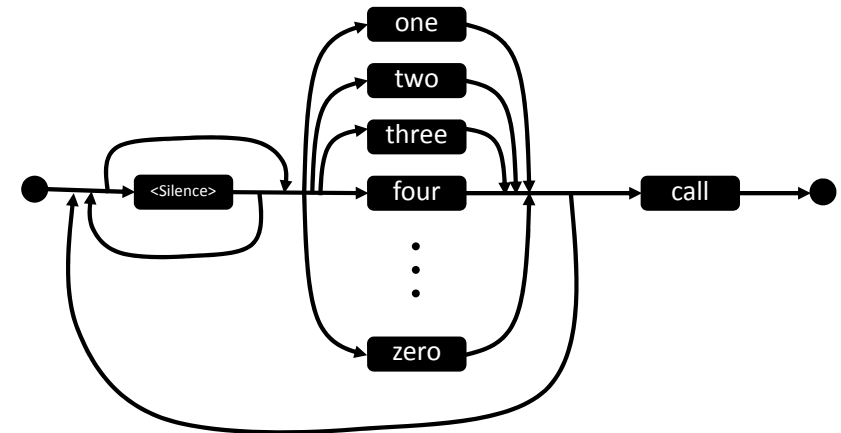
$$P(\text{"buy two tickets for Jandara"}) = P(\text{two})P(\text{Jandara})$$

$$P(\text{"buy one ticket for Die Hard"}) = P(\text{one ticket})P(\text{Die Hard})$$

$$P(\text{"buy one Die Hard"}) = 0$$

Language Model

Finite State Networks



Language Model

n-gram Language Model

$$P(w_1, w_2, \dots, w_L) = \prod_{i=1}^L P(w_i | w_1, w_2, \dots, w_{i-1})$$

$$\approx \prod_{i=1}^L P(w_i | w_{i-(n-1)}, w_{i-(n-2)}, \dots, w_{i-1})$$

Unigram

$$P(w_1, \dots, w_L) \approx \prod_{i=1}^L P(w_i)$$

Bigram

$$P(w_1, \dots, w_L) \approx \prod_{i=1}^L P(w_i | w_{i-1})$$

Trigram

$$P(w_1, \dots, w_L) \approx \prod_{i=1}^L P(w_i | w_{i-1}, w_{i-2})$$

Language Model

$$P(\text{หมีเป็นสัตว์กินเบียร์}) = P(\text{หมี}) \times P(\text{เป็น} | \text{หมี}) \times P(\text{สัตว์} | \text{เป็น, หมี}) \\ \times P(\text{กิน} | \text{สัตว์, เป็น, หมี}) \times P(\text{เบียร์} | \text{กิน, สัตว์, เป็น, หมี})$$

Unigram Approximation

$$P(\text{หมีเป็นสัตว์กินเบียร์}) \approx P(\text{หมี}) \times P(\text{เป็น}) \times P(\text{สัตว์}) \times P(\text{กิน}) \times P(\text{เบียร์})$$

Bigram Approximation

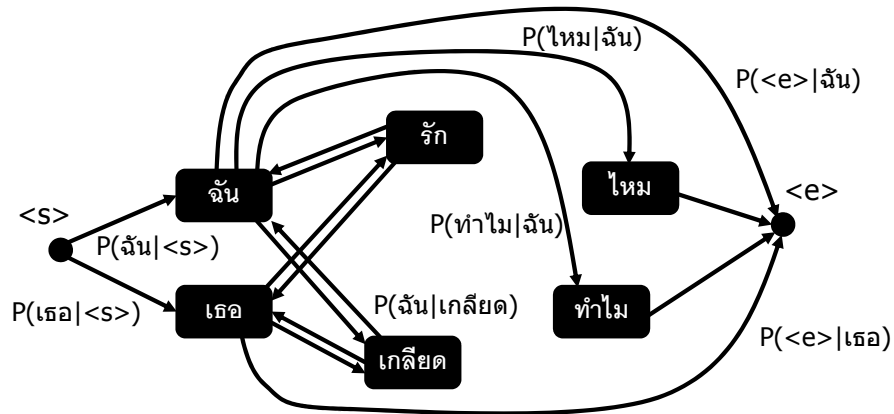
$$P(\text{หมีเป็นสัตว์กินเบียร์}) \approx P(\text{หมี} | \text{<start>}) \times P(\text{เป็น} | \text{หมี}) \times P(\text{สัตว์} | \text{เป็น}) \\ \times P(\text{กิน} | \text{สัตว์}) \times P(\text{เบียร์} | \text{กิน}) \times P(\text{<end>} | \text{เบียร์})$$

Trigram Approximation

$$P(\text{หมีเป็นสัตว์กินเบียร์}) \approx P(\text{หมี} | \text{<start>, <start>}) \times P(\text{เป็น} | \text{หมี, <start>}) \times P(\text{สัตว์} | \text{เป็น, หมี}) \\ \times P(\text{กิน} | \text{สัตว์, เป็น}) \times P(\text{เบียร์} | \text{กิน, สัตว์}) \times P(\text{<end>} | \text{เบียร์, กิน})$$

Language Model

Bigram Approximation



$$P(\text{เรอเกลียดจัน}|X) \propto P(X|\text{เรอเกลียดจัน})P(\text{เรอ}|\text{<S>})P(\text{เกลียด}|\text{เรอ})P(\text{จัน}|\text{เกลียด})P(\text{<e>}|\text{จัน})$$

