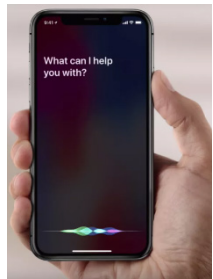


Speech Recognition with Deep Learning Methods

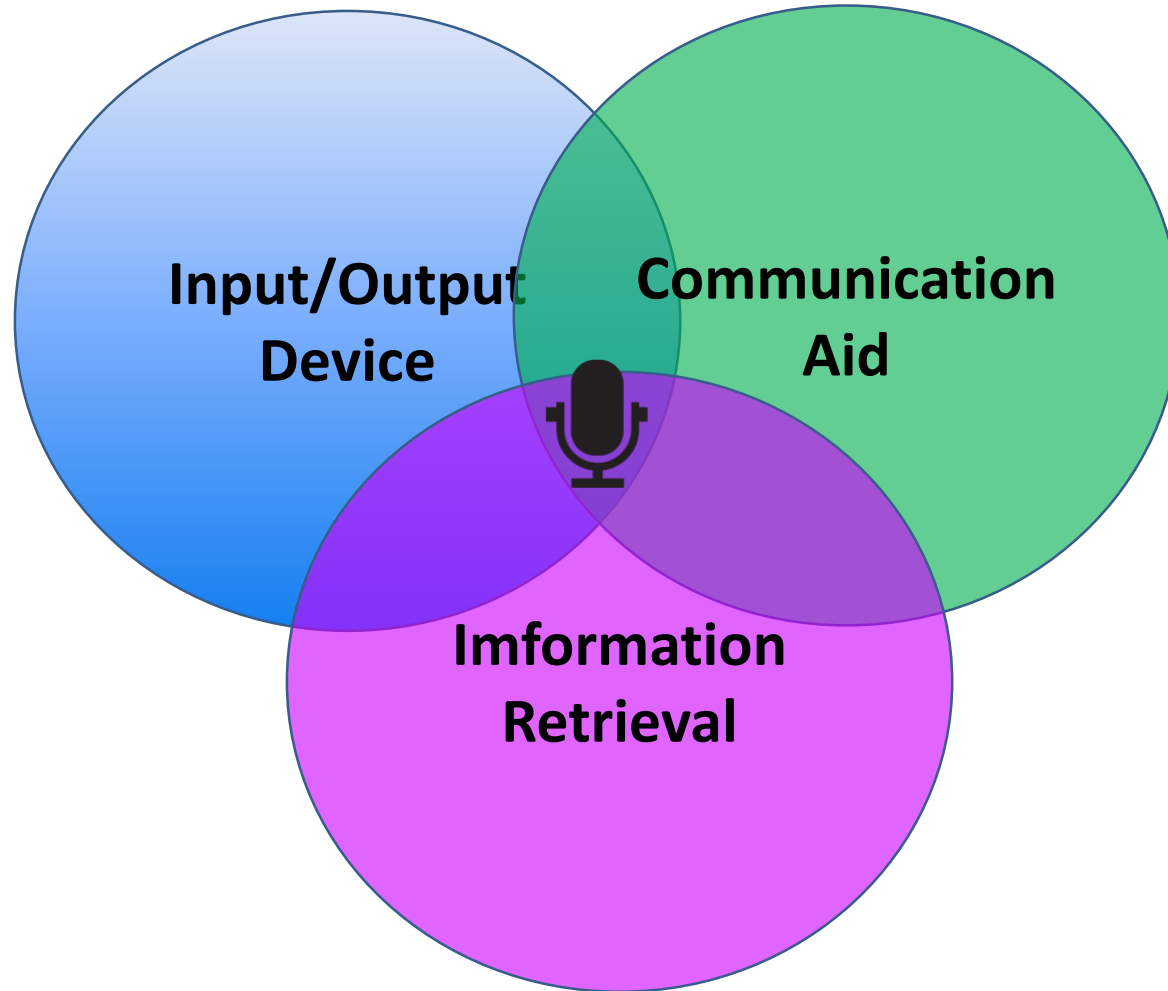
Ming-Han Yang
楊明翰





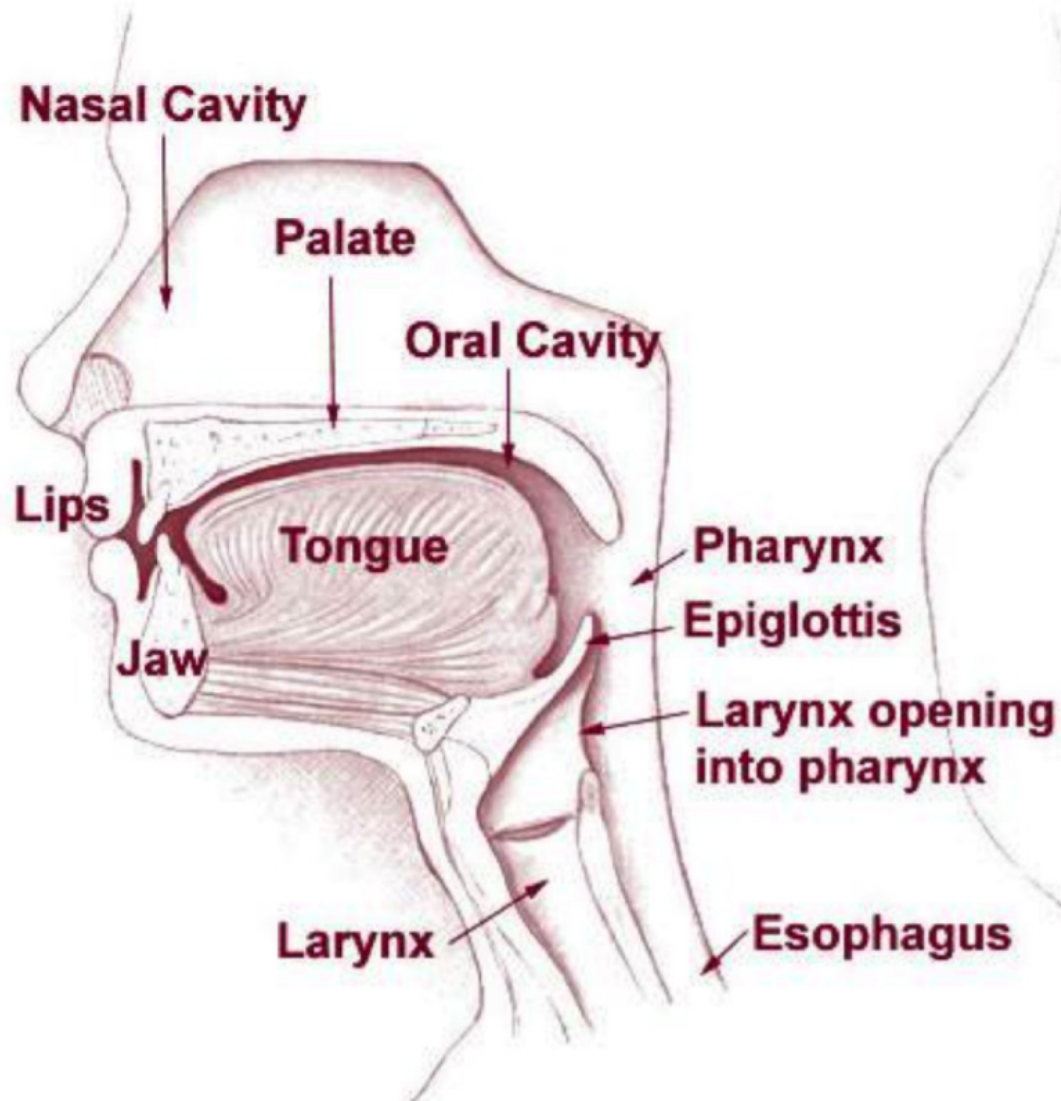


amazon alexa

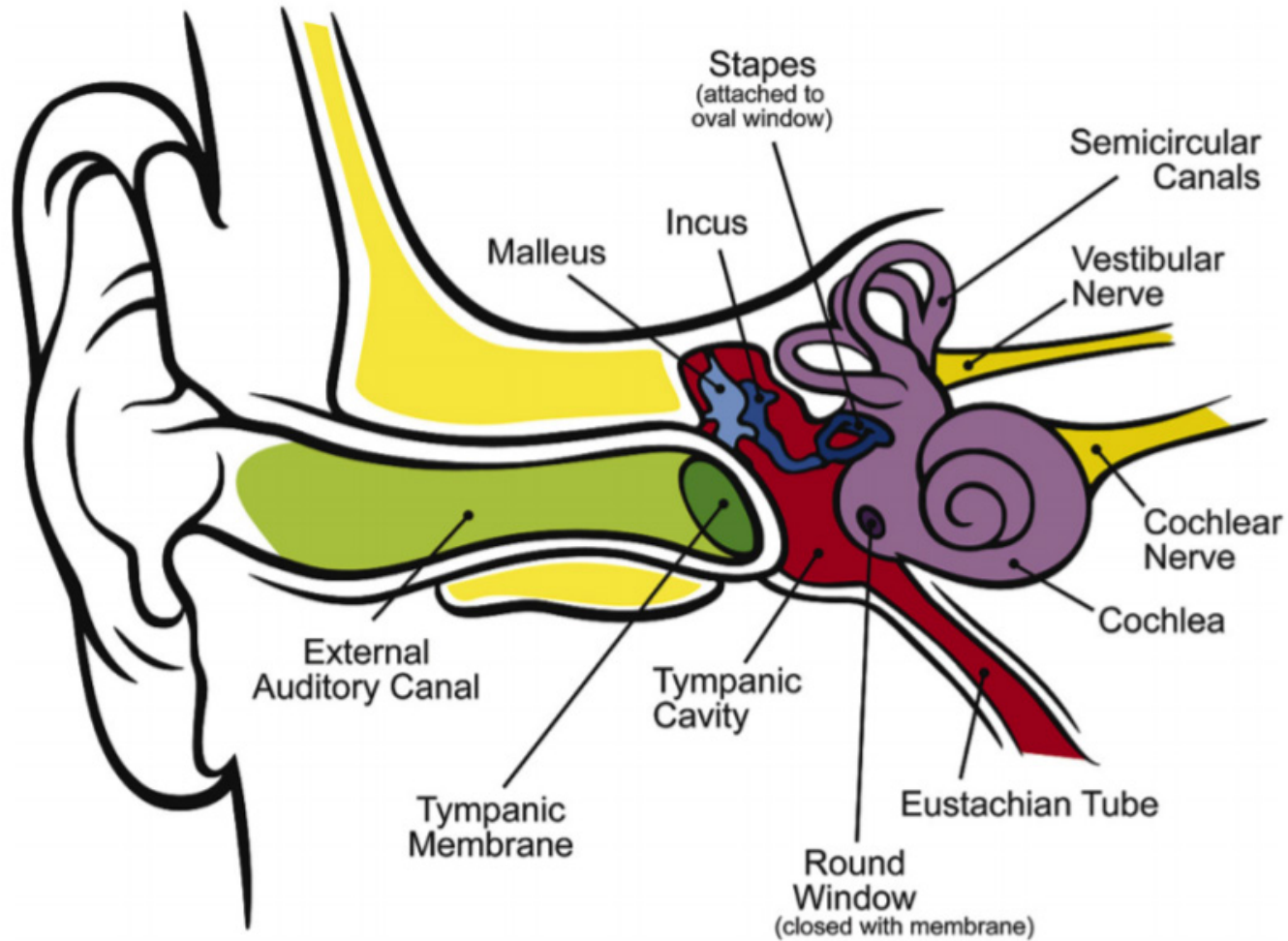


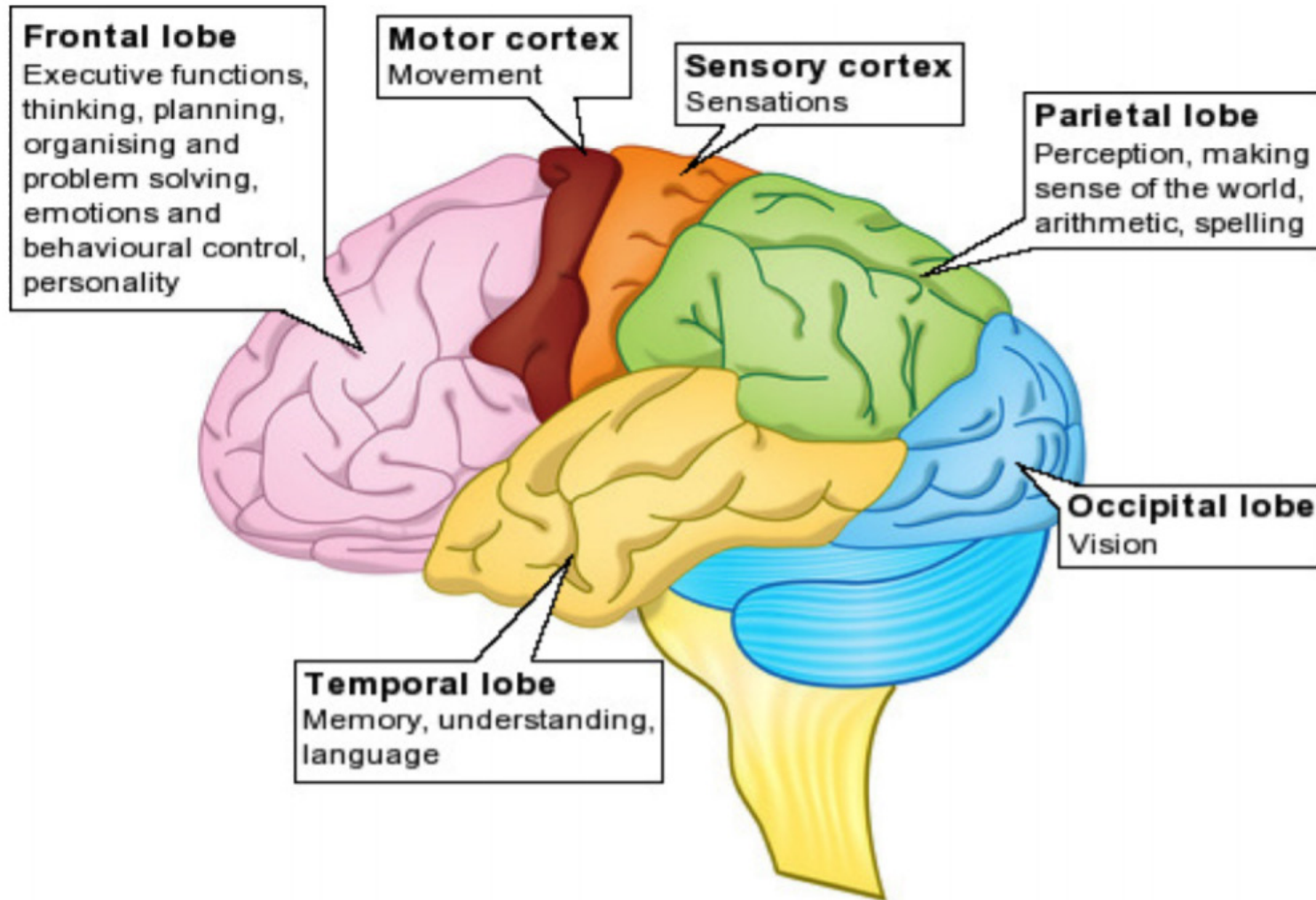

Google Home
Voice-activated speaker





Speech Perception (Recognition)



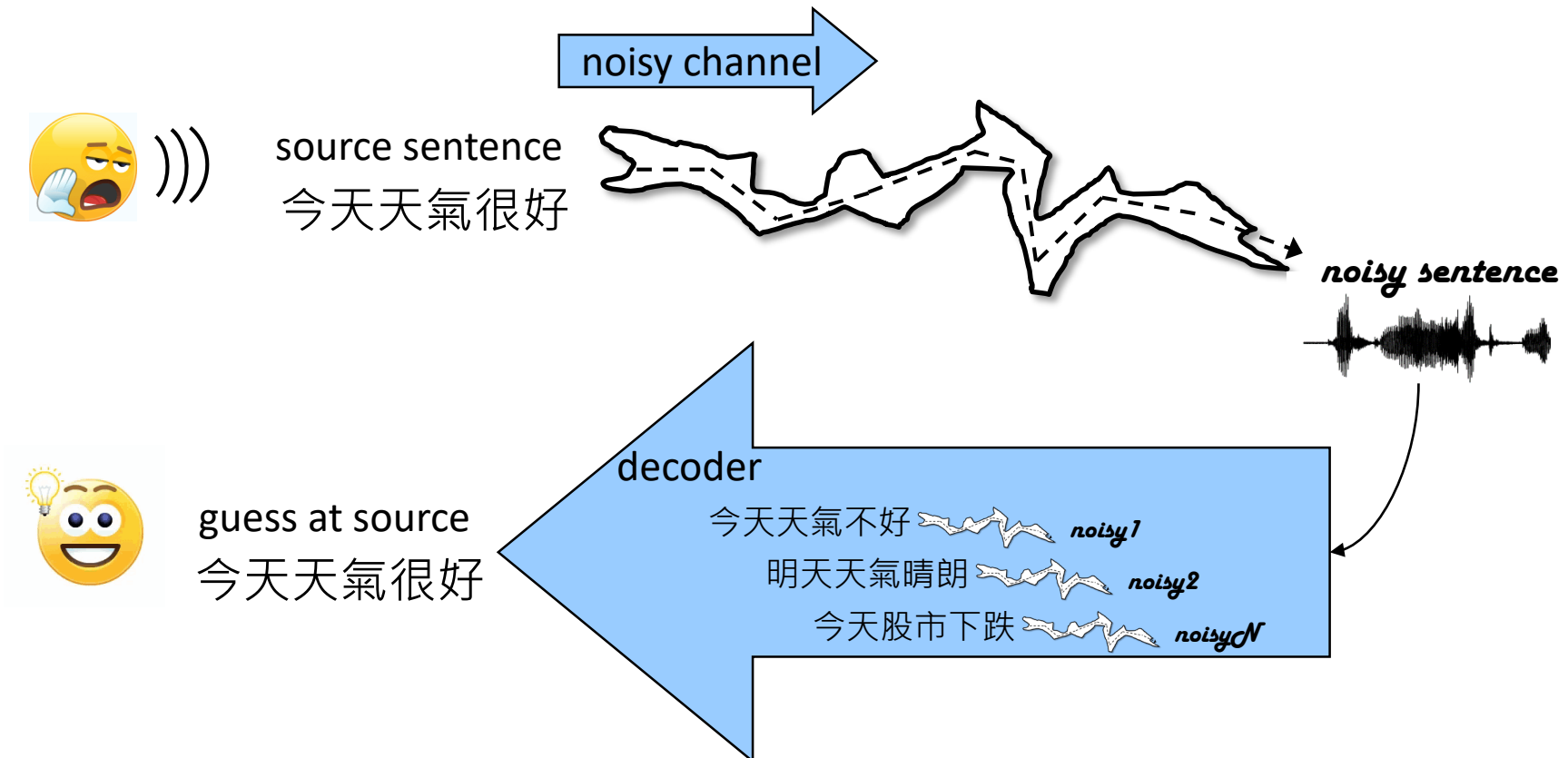


- ◆ Traditional Speech Recognition
- ◆ How to use Deep Learning in acoustic modeling?
- ◆ Why Deep Learning?
- ◆ Speaker Adaptation
- ◆ Multi-task Deep Learning
- ◆ New acoustic features
- ◆ Convolutional Neural Network (CNN)
- ◆ Applications in Acoustic Signal Processing

Traditional Speech Recognition



- ◆ Search through space of all possible sentences.
- ◆ Pick the one that is most probable given the waveform.





Waveform



今天天氣很好

Words



Waveform

Features

今天天氣很好

Words



Waveform

Features

/t/ /i/ /a/ /n/

Phones

今 天 天 氣 很 好

Words



Waveform

Features

/t/-/i/+ /a/

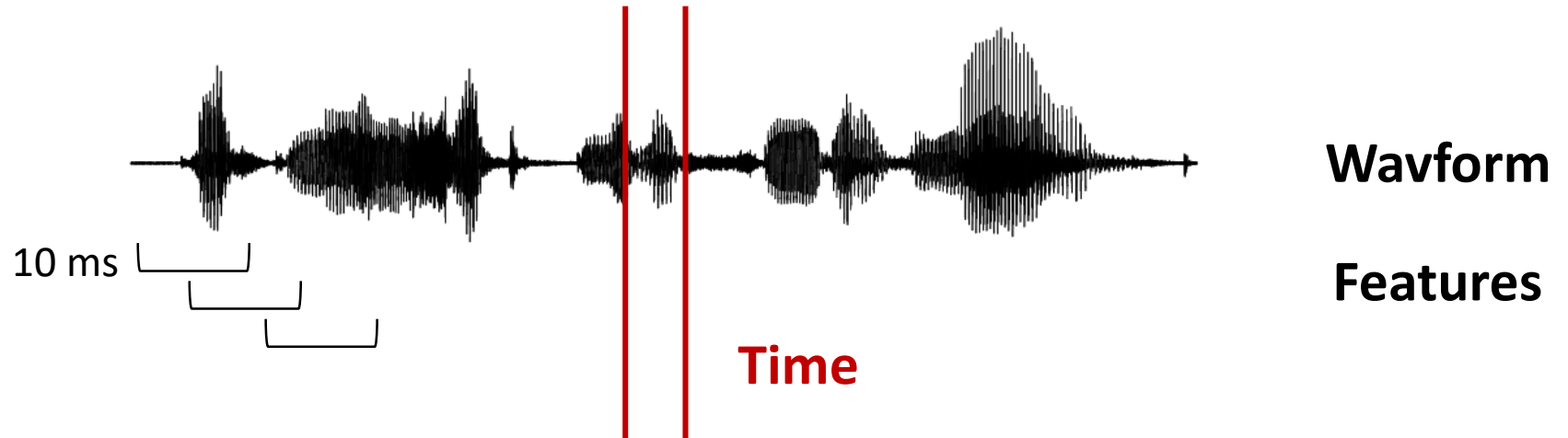
Context-Dependent
Phones

/t/ /i/ /a/ /n/

Phones

今天天氣很好

Words



/t/-/i/+/a/

Context-Dependent
Phones

/t/ /i/ /a/ /n/

Phones

今天天氣很好

Words

◆ Sequence-to-sequence modelling central to speech/language:

● machine translation:

✳ word sequence (discrete) → word sequence (discrete)

– 你好 → Hello

● speech synthesis:

✳ word sequence (discrete) → waveform (continuous)

– 你好 → 

● speech recognition:

✳ waveform (continuous) → word sequence (discrete)

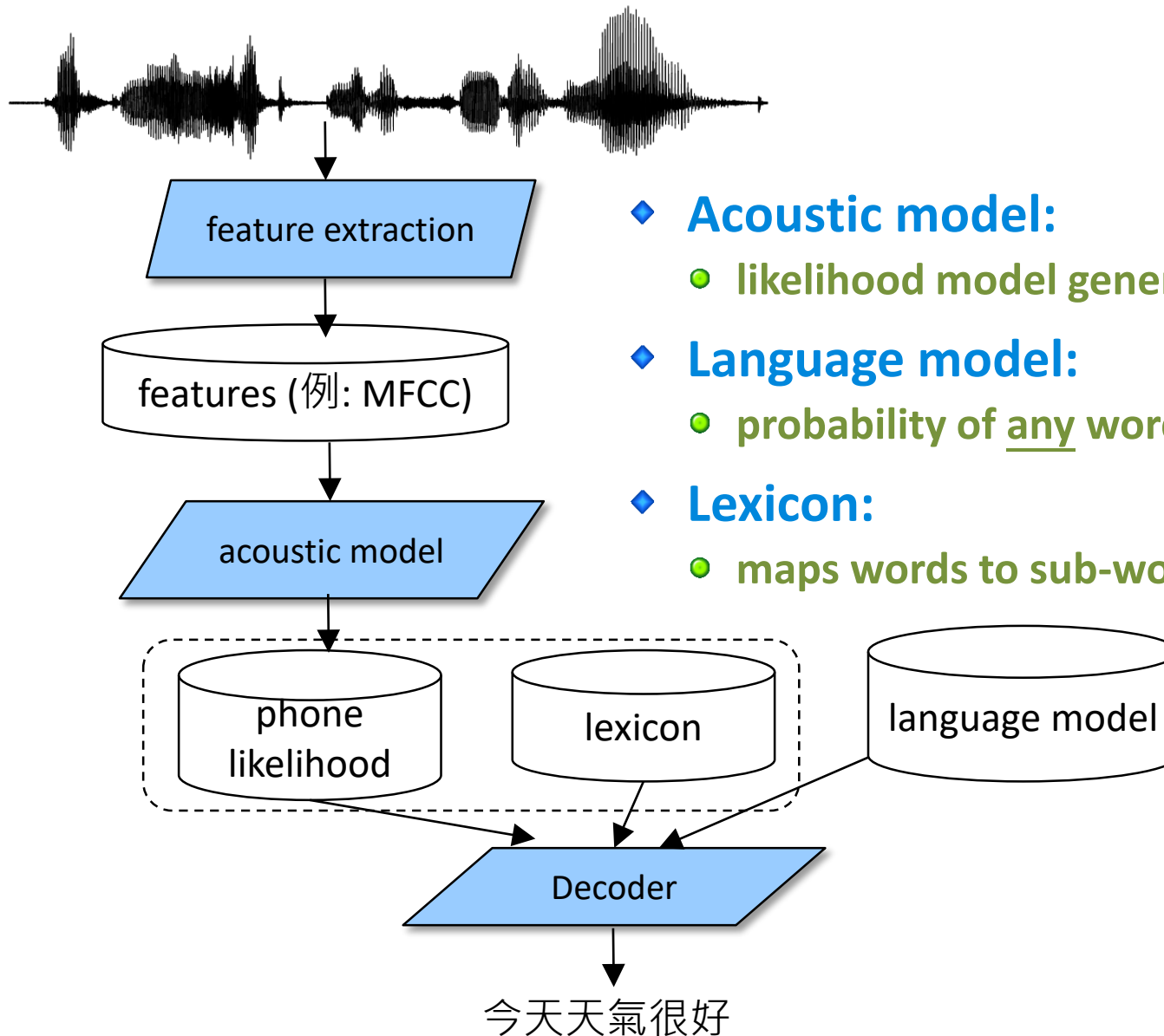
–  → 你好

◆ The sequence lengths on either side can differ

● waveform sampled at 10ms/5ms frame-rate: T -length

● word/token sequences: L -length

● T 遠大於 L



◆ **Acoustic model:**

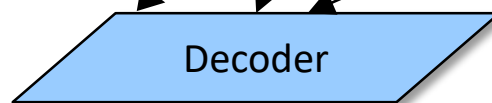
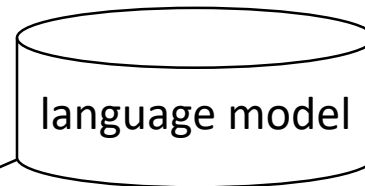
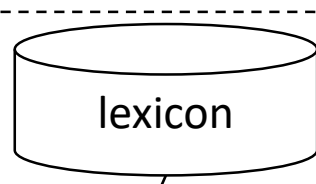
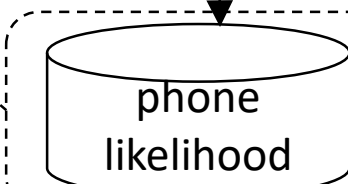
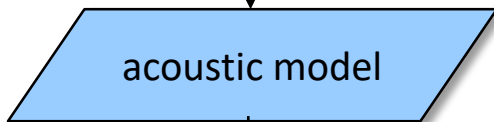
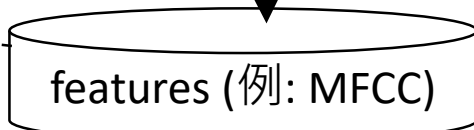
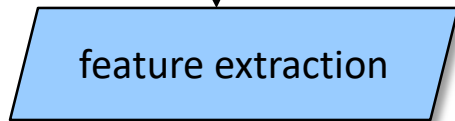
- likelihood model generating observed features

◆ **Language model:**

- probability of any word sequence

◆ **Lexicon:**

- maps words to sub-word units (phones)



\mathbf{W} 今天天氣很好

$$\mathbf{W}^* = \arg \max_{\mathbf{W} \in L} P(\mathbf{W}|\mathbf{O})$$

(貝氏定理) = $\arg \max_{\mathbf{W} \in L} \frac{p(\mathbf{O}|\mathbf{W})P(\mathbf{W})}{p(\mathbf{O})}$

= $\arg \max_{\mathbf{W} \in L} p(\mathbf{O}|\mathbf{W})P(\mathbf{W})$

AM (likelihood)

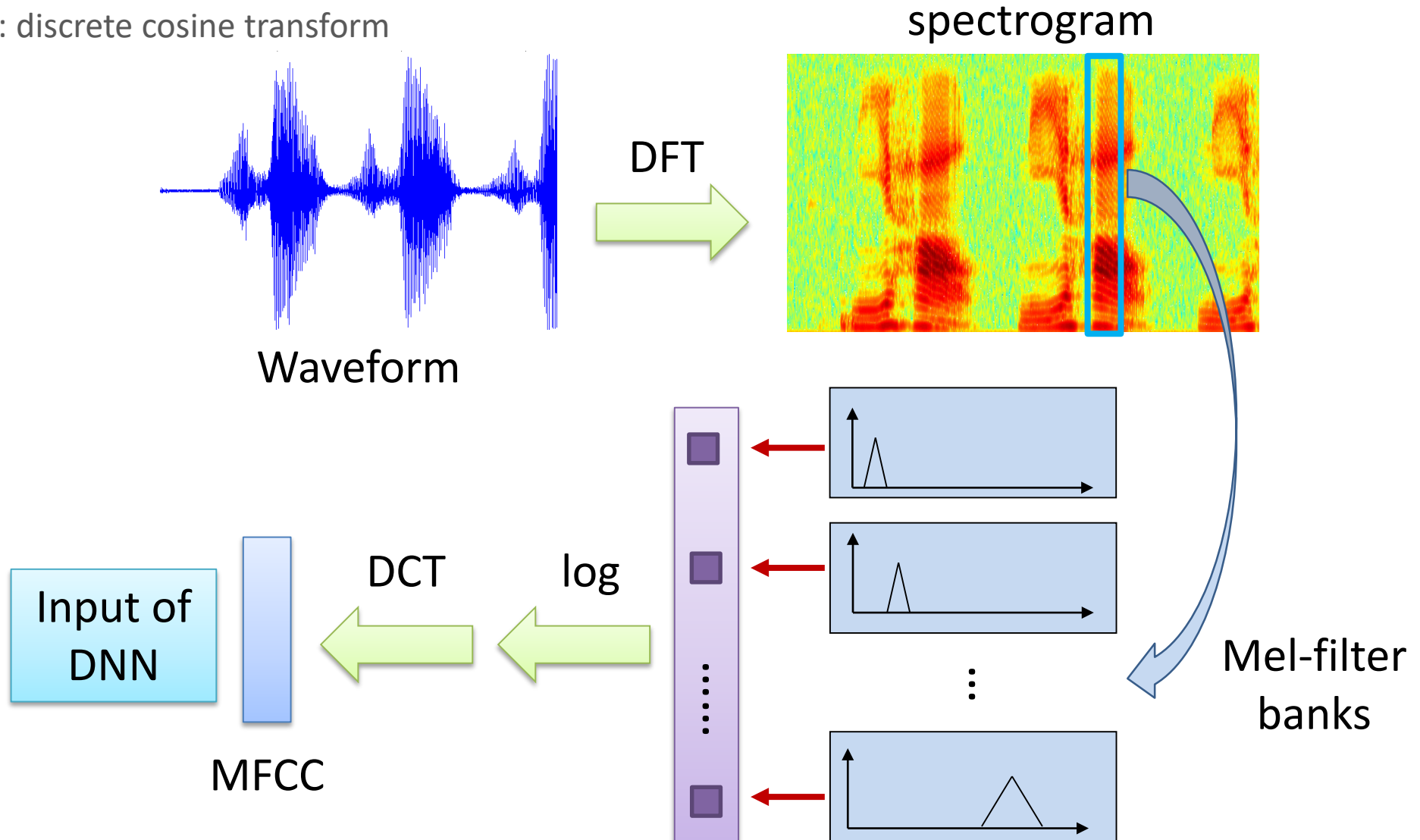
LM (prior)

$P(\mathbf{W})$

$p(\mathbf{O}|\mathbf{W})$

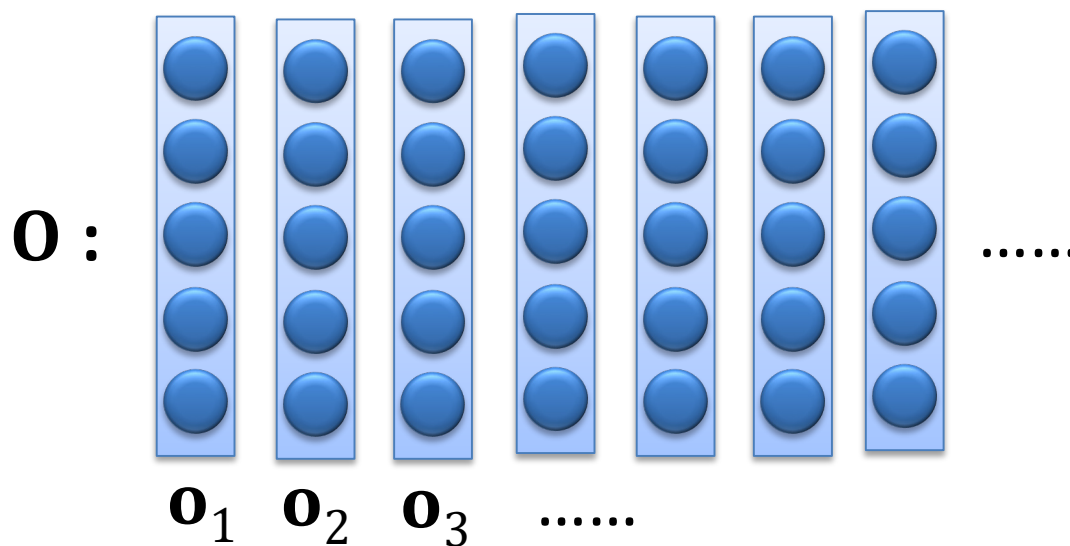
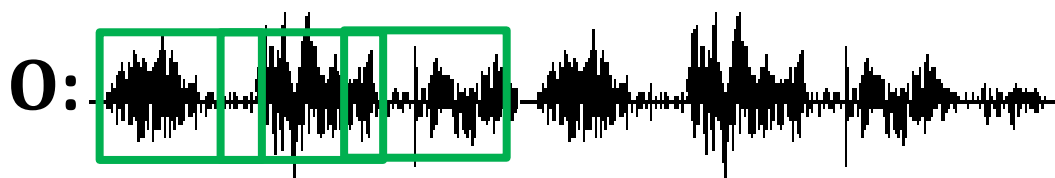
DFT: discrete fourier transform

DCT: discrete cosine transform



◆ Audio is represented by a vector sequence

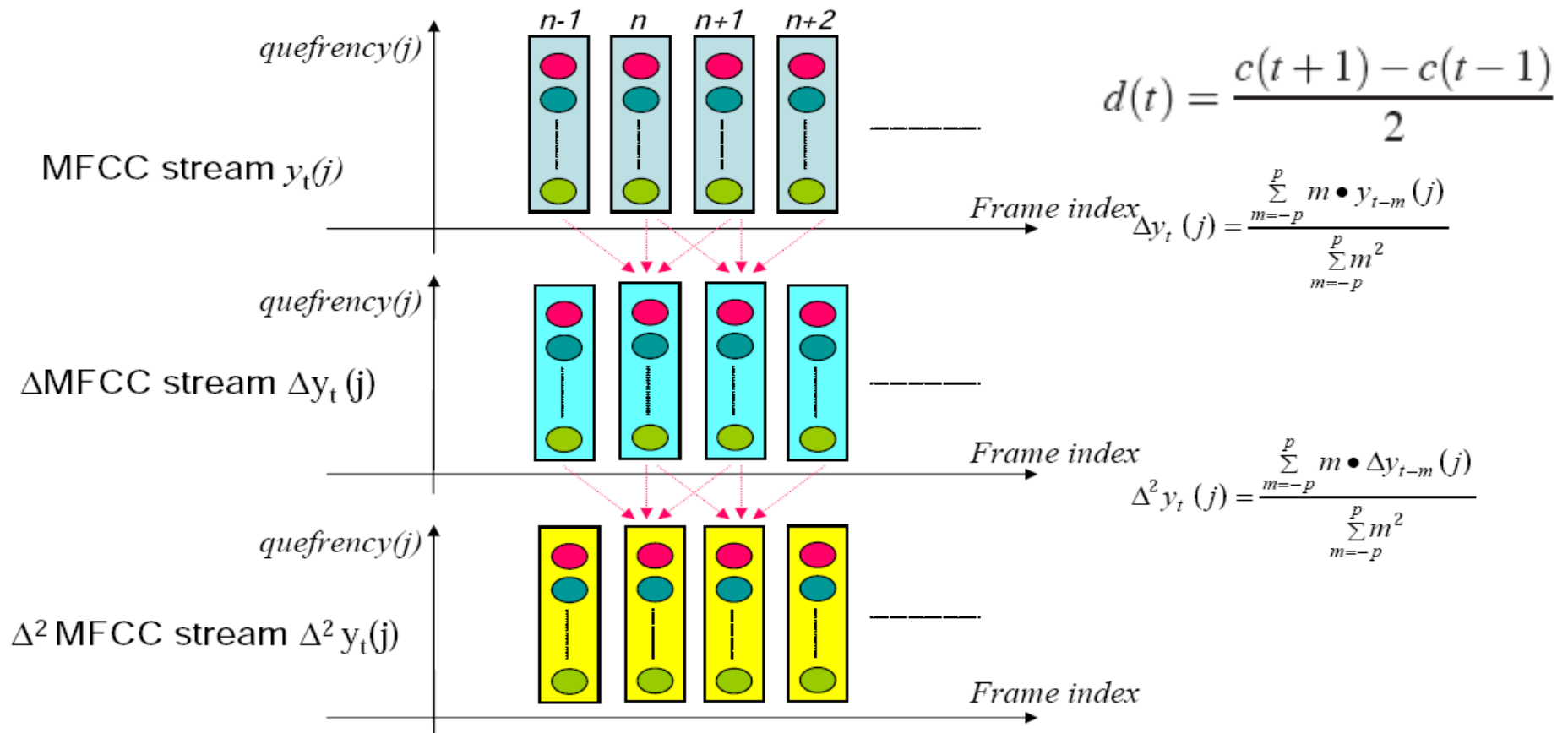
- 聲音 → 音框(frame): $\mathbf{o}_1, \mathbf{o}_2, \mathbf{o}_3, \dots$



每一個frame是
13 dim MFCC

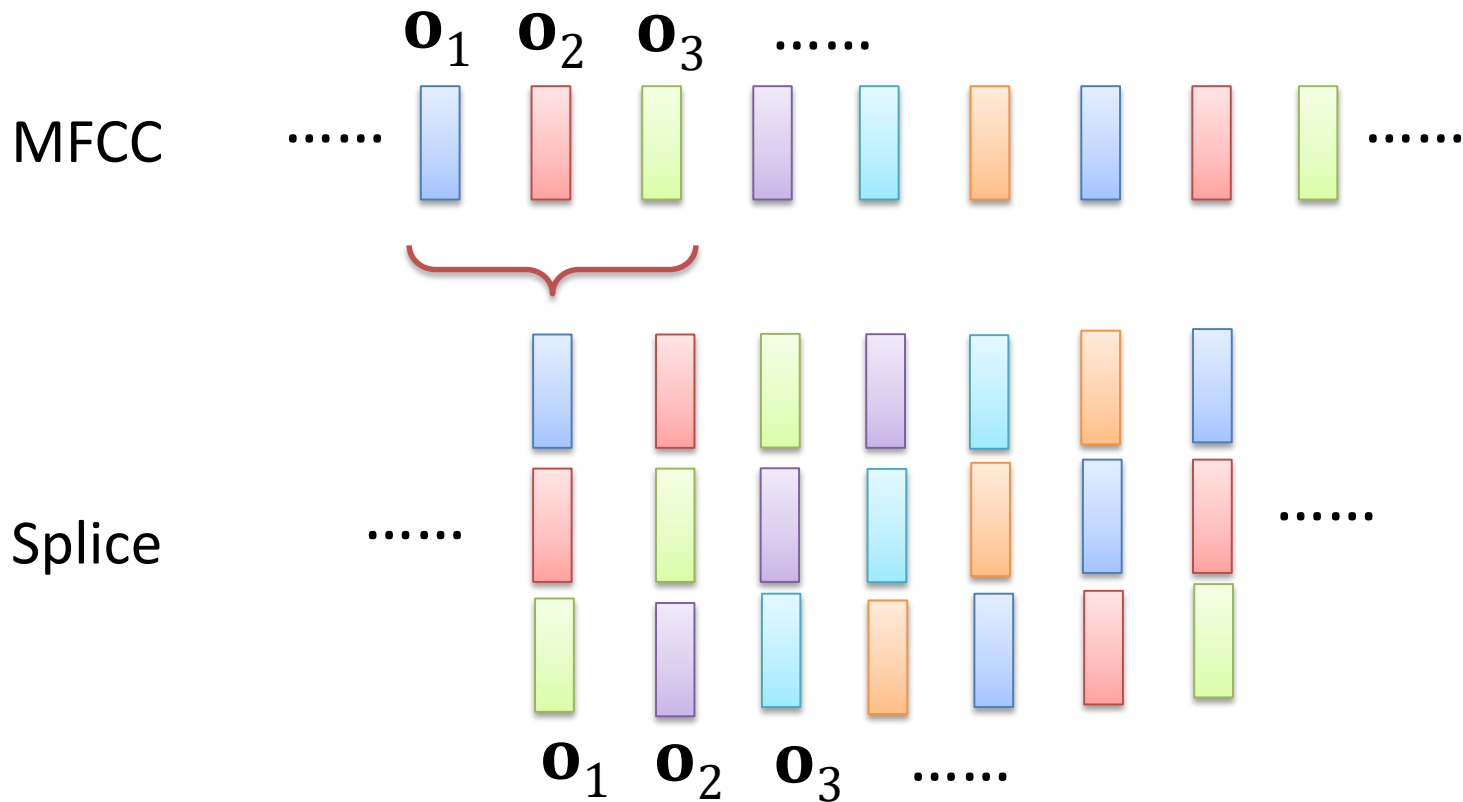
◆ **Derivative: in order to obtain temporal information**

- 可以看成是MFCC的速度與加速度的資訊
- 所以每一個frame會變成 39 dim MFCC



◆ To consider some temporal information

- 把前後的frame跟目前的frame串起來，變成一個新的vector

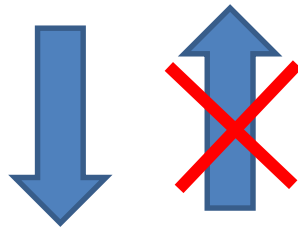


- ◆ **Phoneme: basic unit**
 - **Lexicon: maps words to sub-word units**
- ◆ **Each word corresponds to a sequence of phonemes**

Lexicon

黑面琵鷺	h ei1 m ian4 p i2 l u4
黑馬	h ei1 m a3
黑體	h ei1 t i3
黑髮	h ei1 f a3
黑鯛	h ei1 d iao1
黑鯨	h ei1 j ing1
context	k aa n t eh k s t
everyday	eh v r iy d ey
include	ih n k l uw d

Word what do you think



Different words can correspond to the same phonemes

Lexicon hh w aa t d uw y uw th ih ng k

- ◆ 利用語言模型, 可以限制Acoustic Model找到的Phone Sequence, 組合起來要長的像人話
- ◆ 常用的LM: N-grams

<s> the cat sat on the mat </s>

<s> the cat sat on the mat **</s>**

<s> the **cat** sat on the mat **</s>**

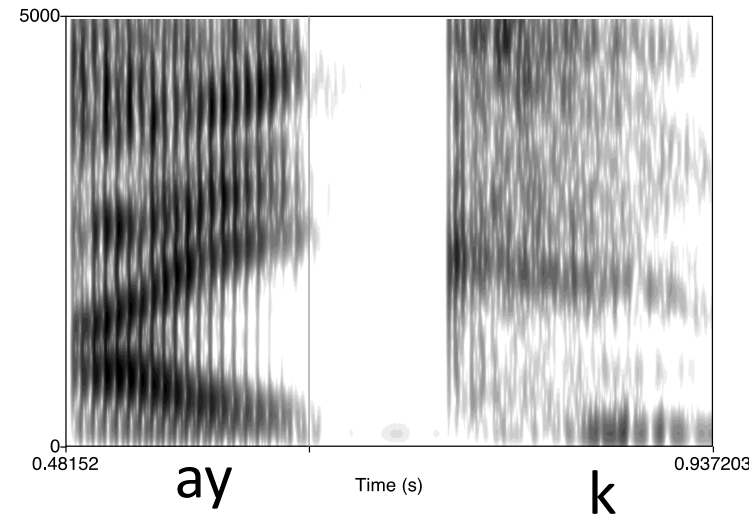
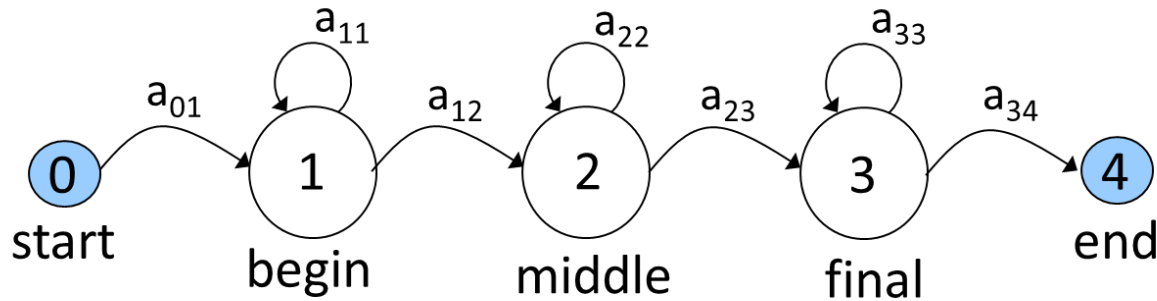
<s> the **cat sat** on the mat **</s>**

<s> the **cat sat on** the mat **</s>**

<s> the cat **sat on the** mat **</s>**

<s> the cat sat **on the mat** **</s>**

<s> the cat sat on **the mat** **</s>**



- ◆ Important sequence model: **Hidden Markov Model (HMM)**
- ◆ HMMs standard model for many year (1970s-2010s)
 - each (context-dependent) phone modelled by an HMM
 - typically 3-emitting state topology, **left-to-right HMM**
 - non-emitting (end) states used for “gluing” models together

- ◆ Each phoneme correspond to a sequence of states

what do you think



Phone:

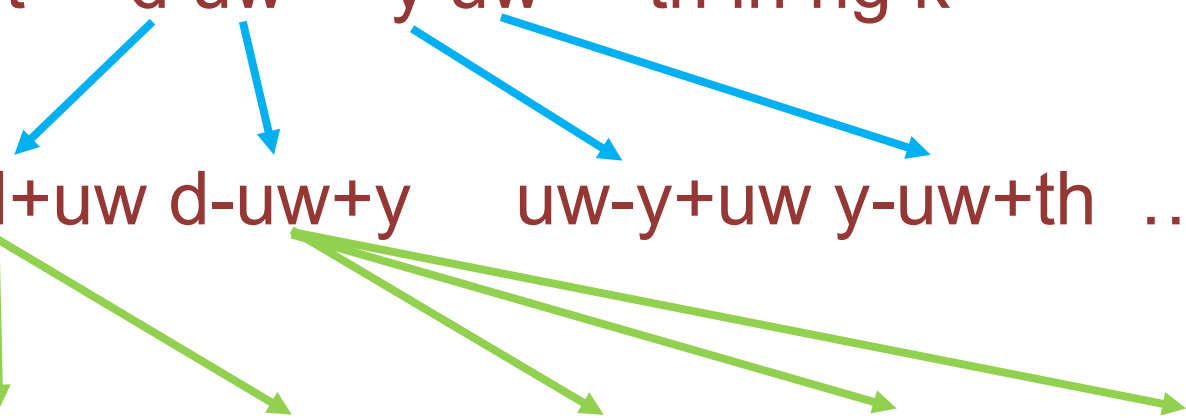
hh w aa t d uw y uw th ih ng k

Tri-phone:

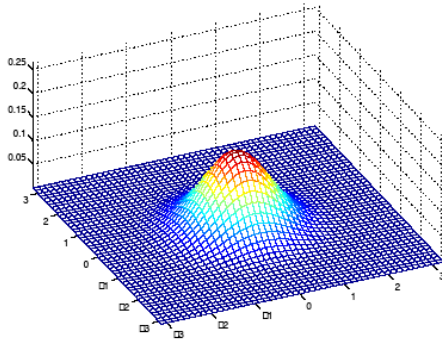
..... t-d+uw d-uw+y uw-y+uw y-uw+th

State:

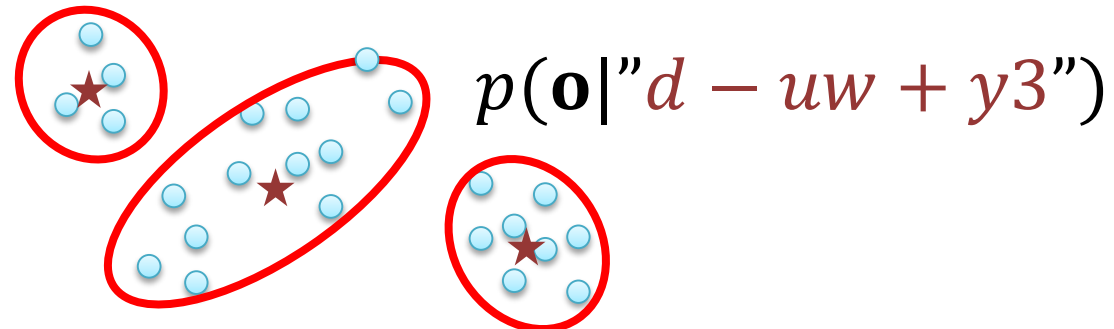
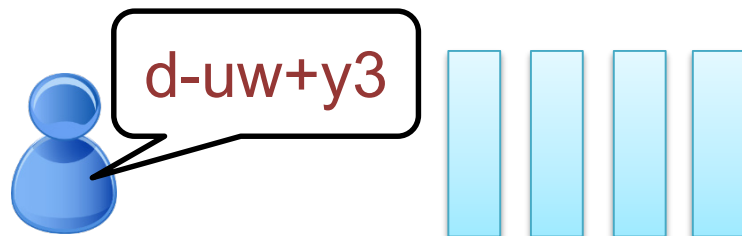
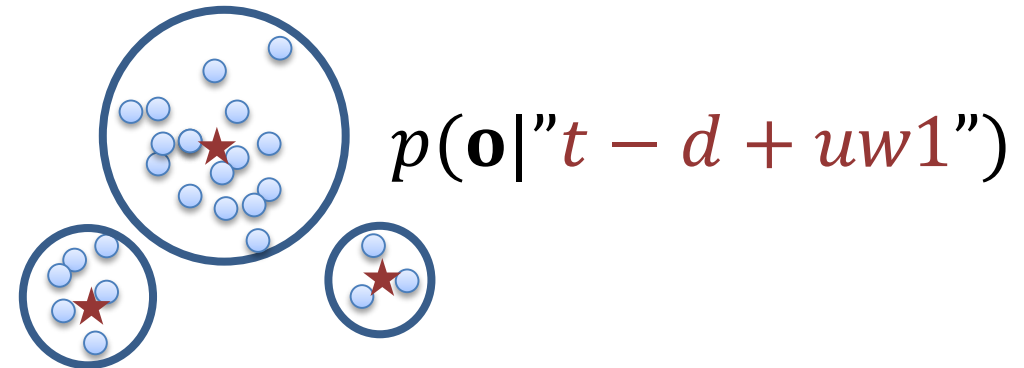
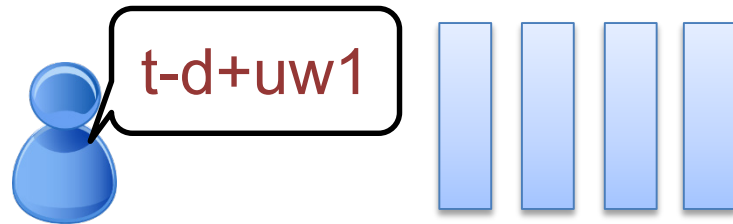
t-d+uw1 t-d+uw2 t-d+uw3 d-uw+y1 d-uw+y2 d-uw+y3



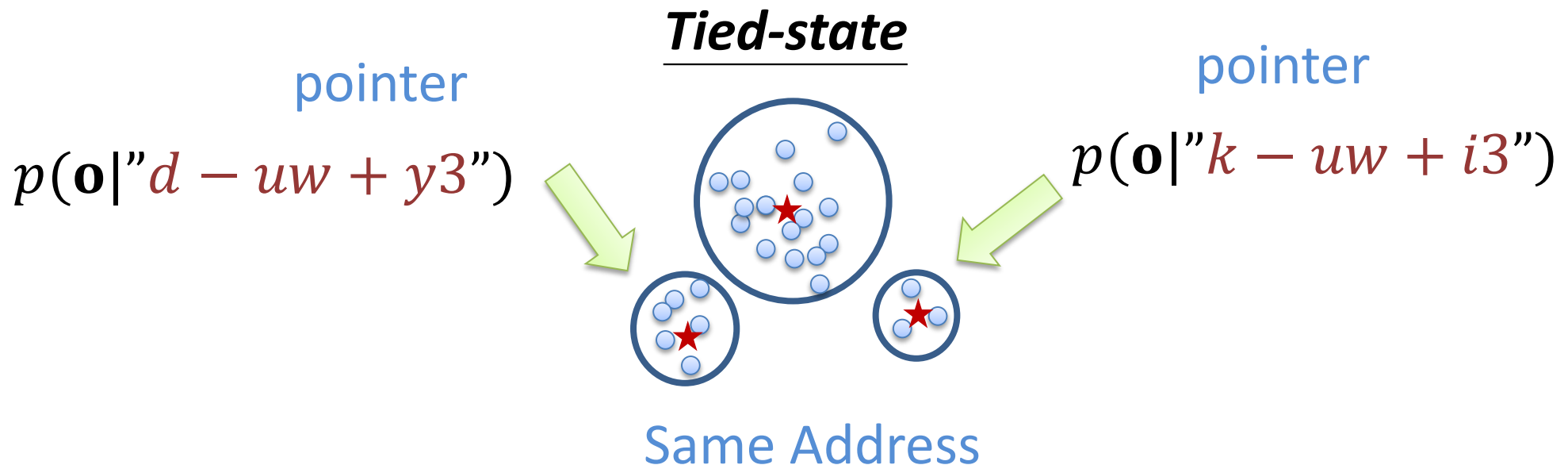
- ◆ Each state has a stationary distribution for acoustic features



Gaussian Mixture Model (GMM)



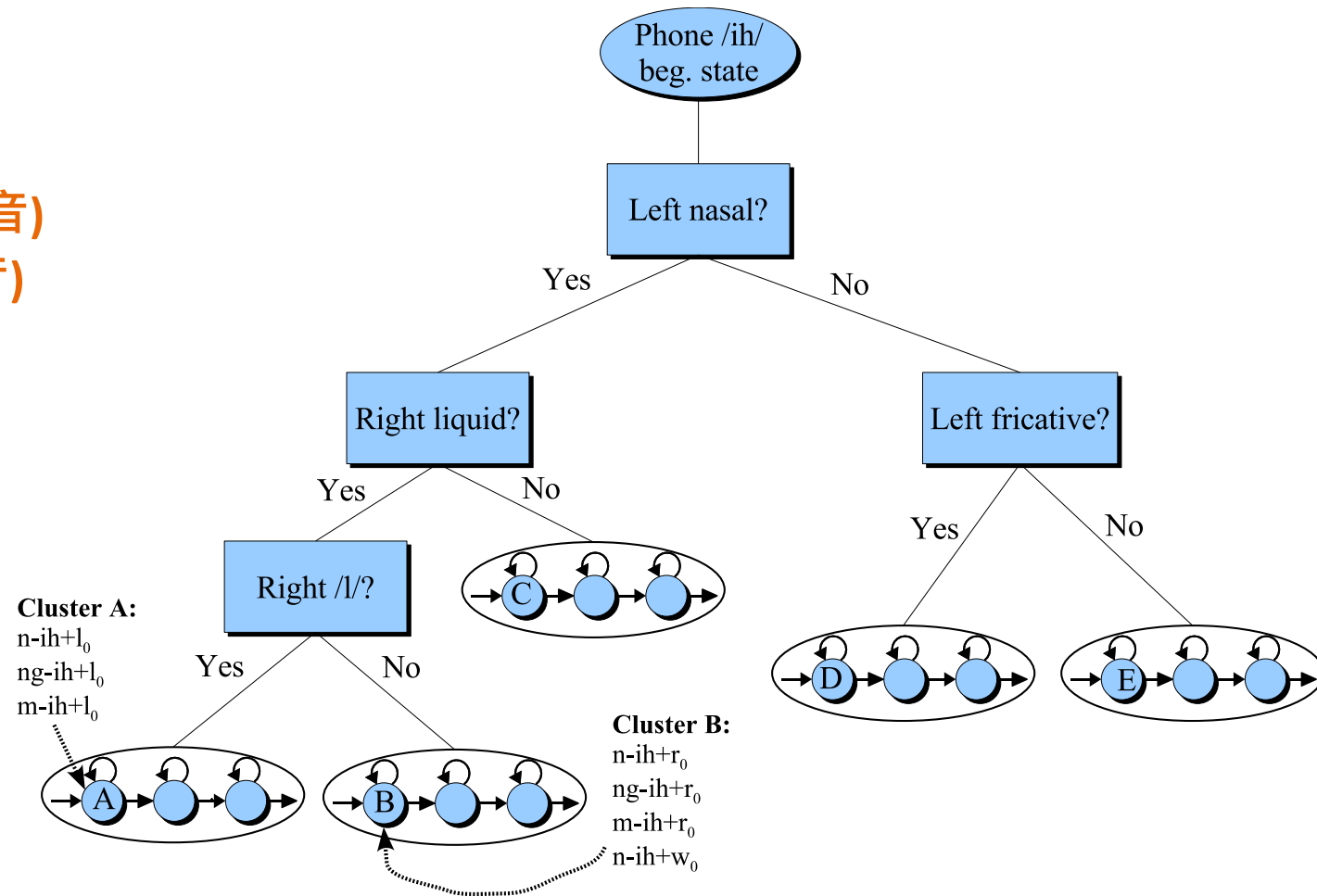
- ◆ 我們可以暴力組合出所有triphone :
 - 假設有50個phone, 總共就會有 $50 \times 50 \times 50 = 125,000$ 個tri-phone
 - 實際用不到那麼多
- ◆ States which are clustered together will share their Gaussians
 - Decision-Tree based clustering of triphone states



◆ How do we decide which triphones to cluster together?

- 可以用語言學家定義的broad phonetic classes

- * Stop
- * Nasal (鼻音)
- * Liquid (流音)
- * Fricative (摩擦音)
- * Sibilant (齒擦音)
- * Vowel (母音)
- * Lateral (邊音)
- * ...

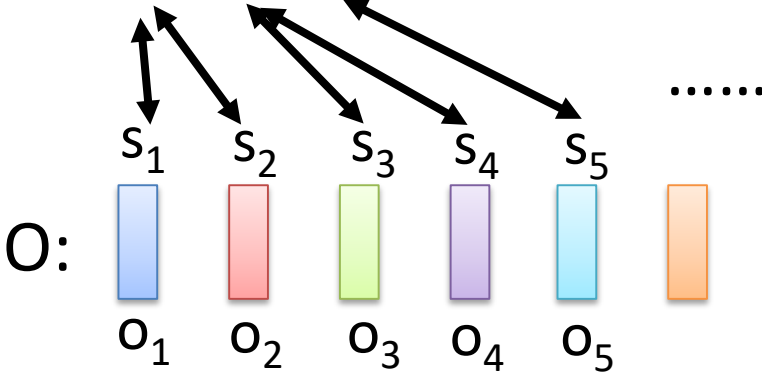
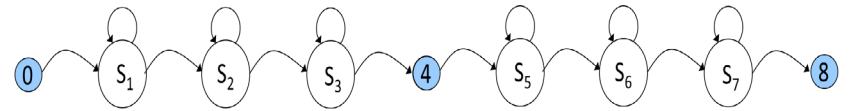


$$W^* = \arg \max_{W \in L} p(\mathbf{O}|W) P(W)$$

W: what do you think?

$$p(\mathbf{O}|W) = P(\mathbf{O}|S)$$

S: a b c d e



Assume we also know the alignment $s_1 \cdots s_T$

$$p(\mathbf{O}|S) = \prod_{t=1}^T \underbrace{p(s_t|s_{t-1})}_{\text{transition}} \underbrace{p(\mathbf{o}_t|s_t)}_{\text{emission}}$$

$$W^* = \arg \max_{W \in L} p(\mathbf{O} | \mathbf{W}) P(\mathbf{W})$$

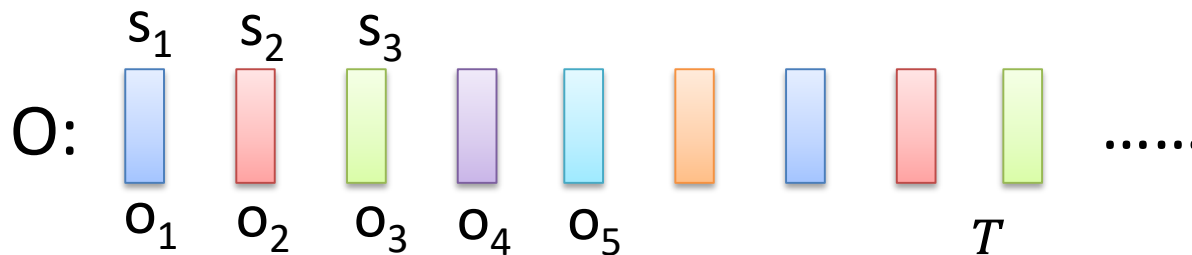
W: what do you think?

$$p(\mathbf{O} | \mathbf{W}) = P(\mathbf{O} | S)$$



S: a b c d e

Actually, we don't know the alignment 😞 → Viterbi algorithm



$$p(\mathbf{O} | S) \approx \max_{S_1 \dots S_T} \prod_{t=1}^T p(s_t | s_{t-1}) p(\mathbf{o}_t | s_t)$$

- Viterbi algorithm is used to find **the most probable path** through a probabilistically scored time/state lattice

- ◆ How to evaluate the word string output by a speech recognizer?

- ◆ **Word Error Rate!**

$$WER = \frac{100\% \times (\textit{Insertions} + \textit{Substitutions} + \textit{Deletions})}{\textit{Total Word in Correct Transcript}}$$

- Insertion: 多一個字, deletion: 少一個字, substitution: 字錯
- Example:

REF:	portable	****	PHONE	UPSTAIRS	last	night	so
HYP:	portable	FORM	OF	STORES	last	night	so
Eval		I	S	S			

* $WER = 100\% \times (1 + 2 + 0)/6 = 50\%$

- ◆ **WER越低，表示辨識器的效果越好！**

◆ Feature Extraction:

- 39 “MFCC” features

◆ Acoustic Model:

- Gaussians for computing phone likelihood $p(O|S)$

◆ Lexicon/Pronunciation Model

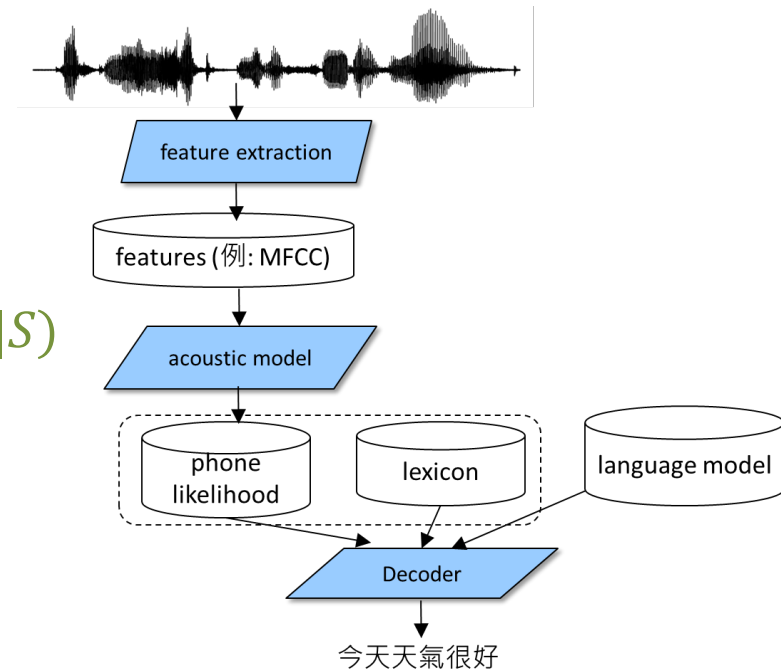
- HMM: what phones can follow each other

◆ Language Model

- N-grams for computing $p(w_i|w_{i-1})$

◆ Decoder

- **Viterbi algorithm**: dynamic programming for combining all these to get word sequence from speech



How to use Deep Learning?

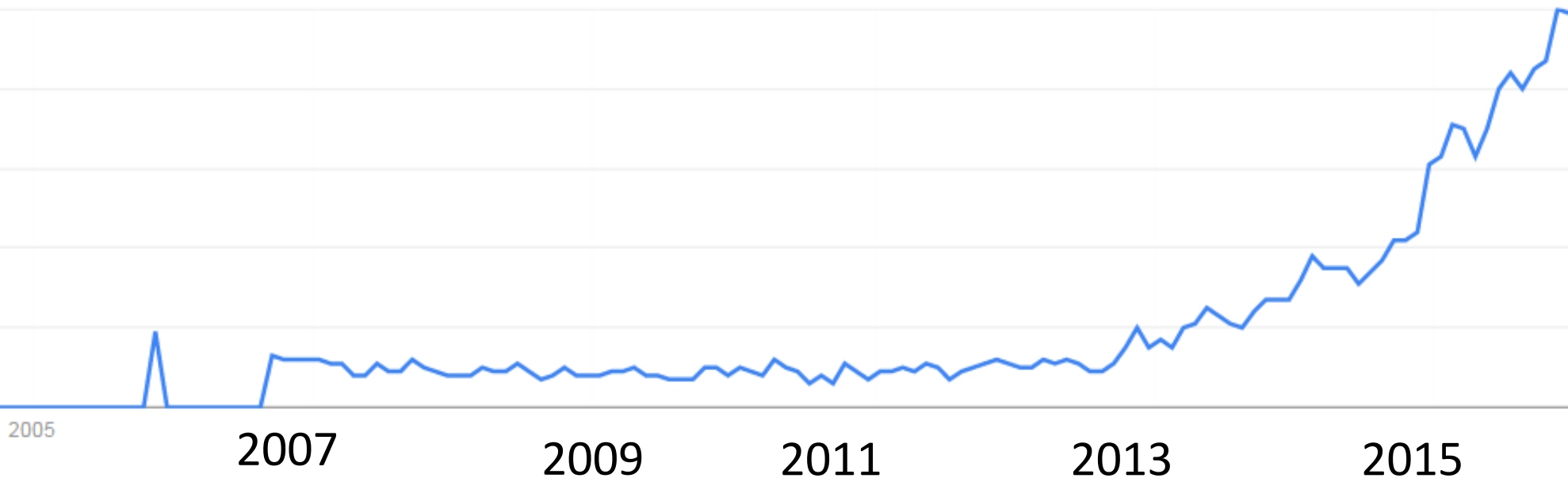


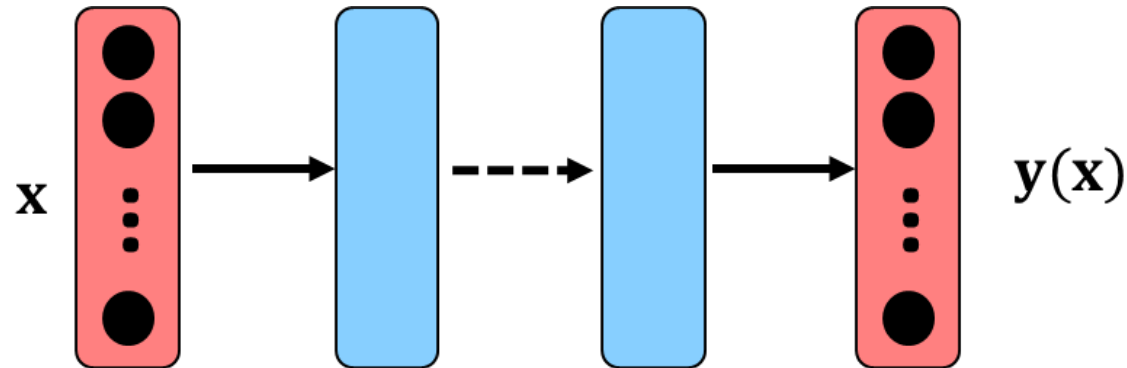
- ◆ **Deep learning attracts lots of attention**

- Deep learning obtains many exciting results

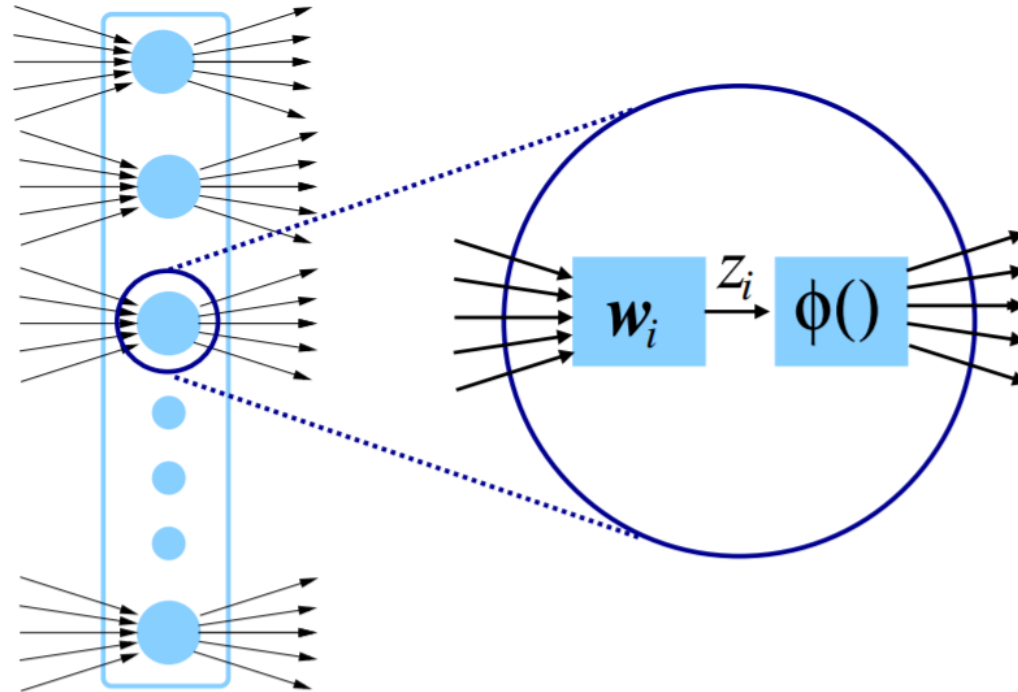
- ◆ **From Wikipedia:**

- 深度學習 (deep learning) 是機器學習的分支，是一種試圖使用包含複雜結構或由多重非線性變換構成的多個處理層對資料進行高層抽象的演算法。



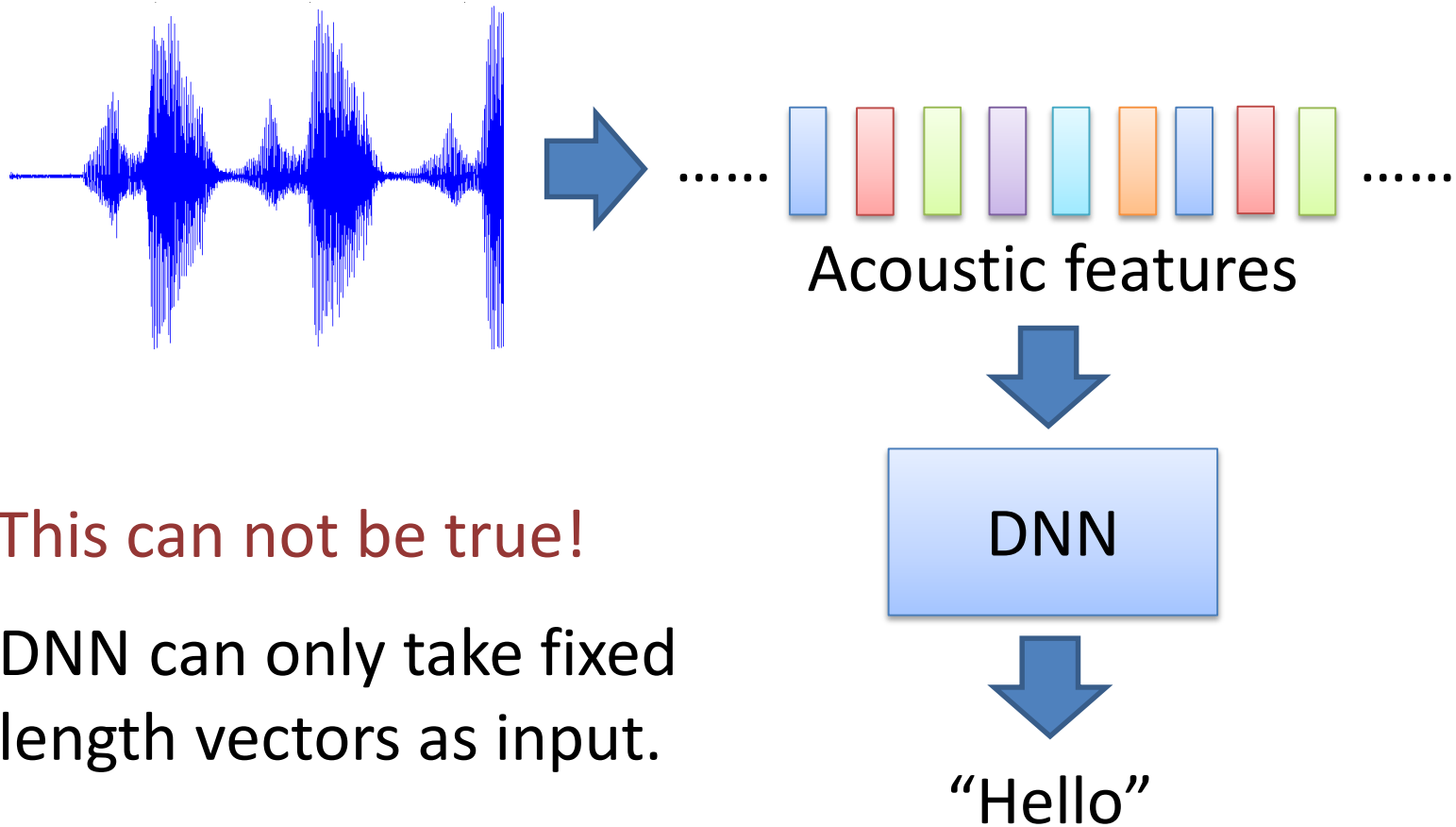


- ◆ **General mapping process from input x to output $y(x)$**
 - deep refers to number of hidden layers
- ◆ **Output from the previous layer connected to following layer:**
 - $x^{(k)}$ is the input to layer k
 - $x^{(k+1)} = y^{(k)}$ the output from layer k



◆ General form for layer k:

- $y_i^{(k)} = \phi(\mathbf{w}_i \mathbf{x}^{(k)} + \mathbf{b}_i) = \phi(z_i^{(k)})$
 - * \mathbf{w} : weight matrix
 - * \mathbf{b} : bias vector
 - * ϕ : activation function
 - Sigmoid, RELU, ...etc



This can not be true!

DNN can only take fixed length vectors as input.

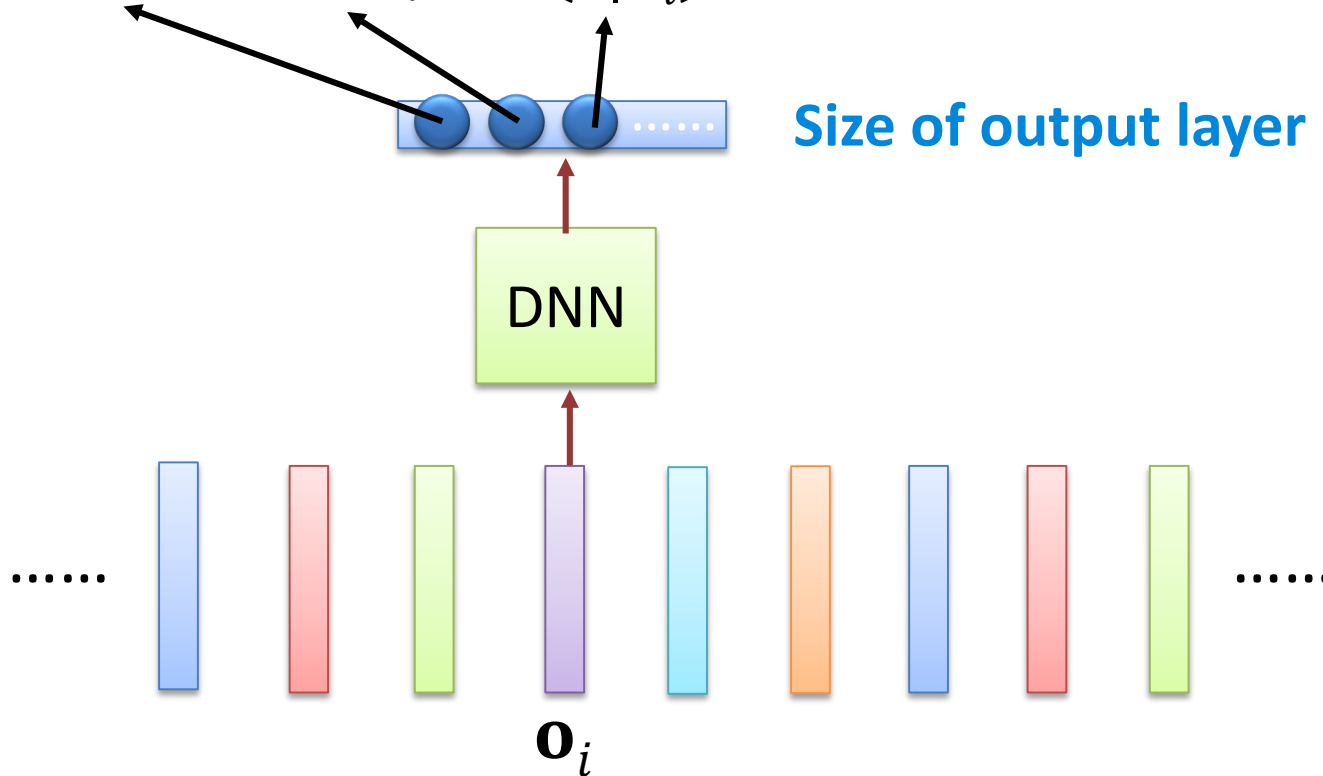
◆ **DNN output:**

- Probability of each state

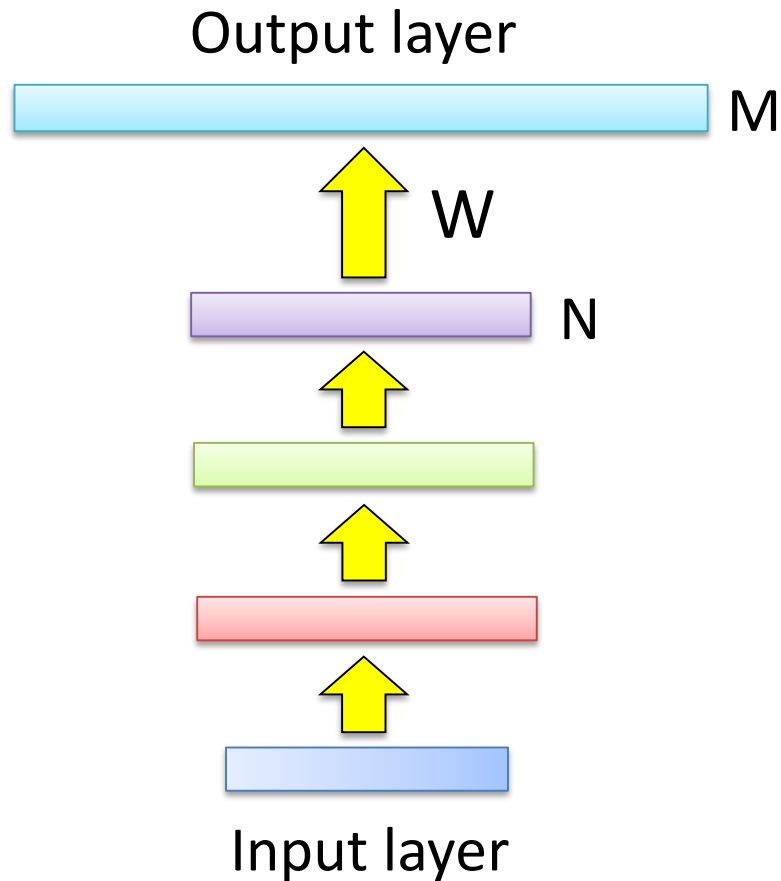
◆ **DNN input:**

- One acoustic feature

$P(a|\mathbf{o}_i)$ $P(b|\mathbf{o}_i)$ $P(c|\mathbf{o}_i)$

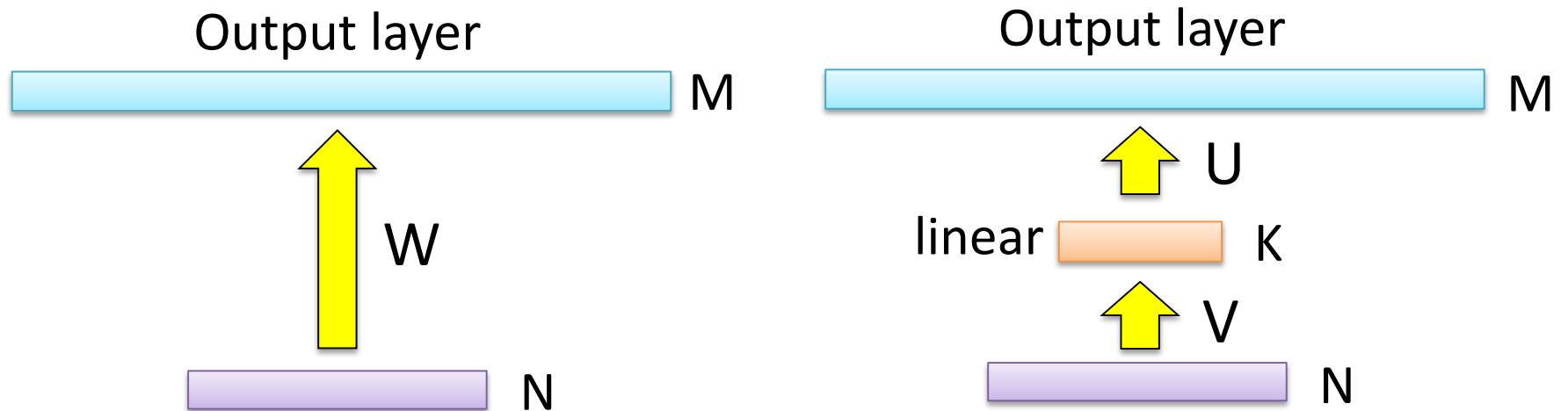
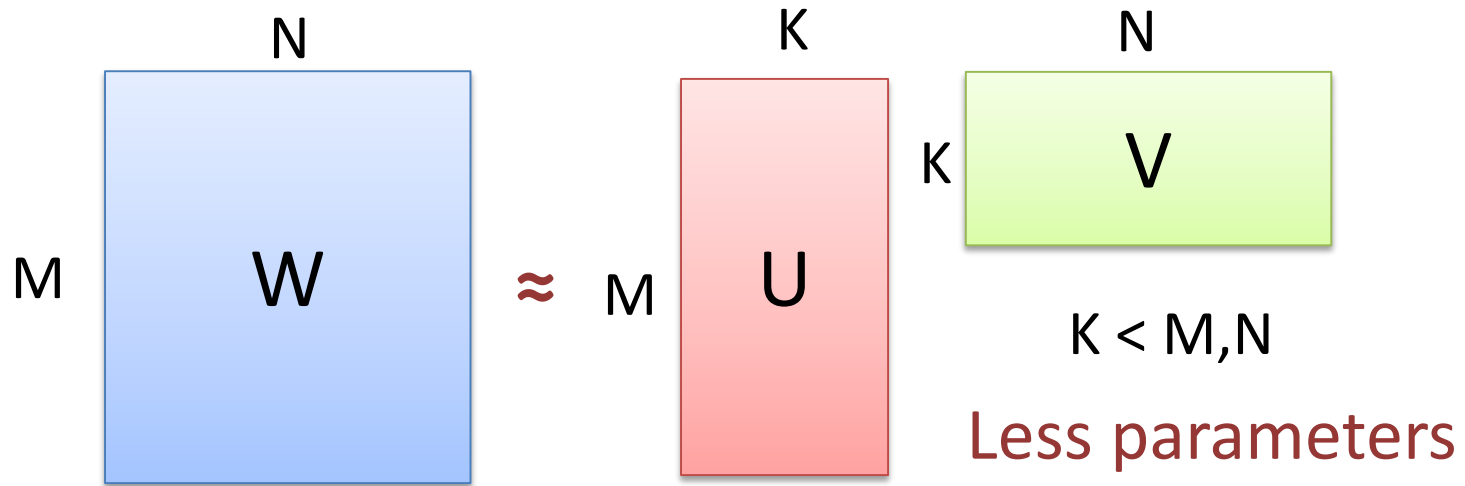


Size of output layer = Num of states



$W: M \times N$

- ◆ **N is the size of the last hidden layer**
- ◆ **M is the size of output layer**
 - Number of states
- ◆ **M can be large if the outputs are the states of tri-phone**
 - e.g. 3000~5000



- ◆ There are three ways to use DNN for acoustic modeling

- Way 1. Tandem
- Way 2. DNN-HMM hybrid
- Way 3. End-to-end

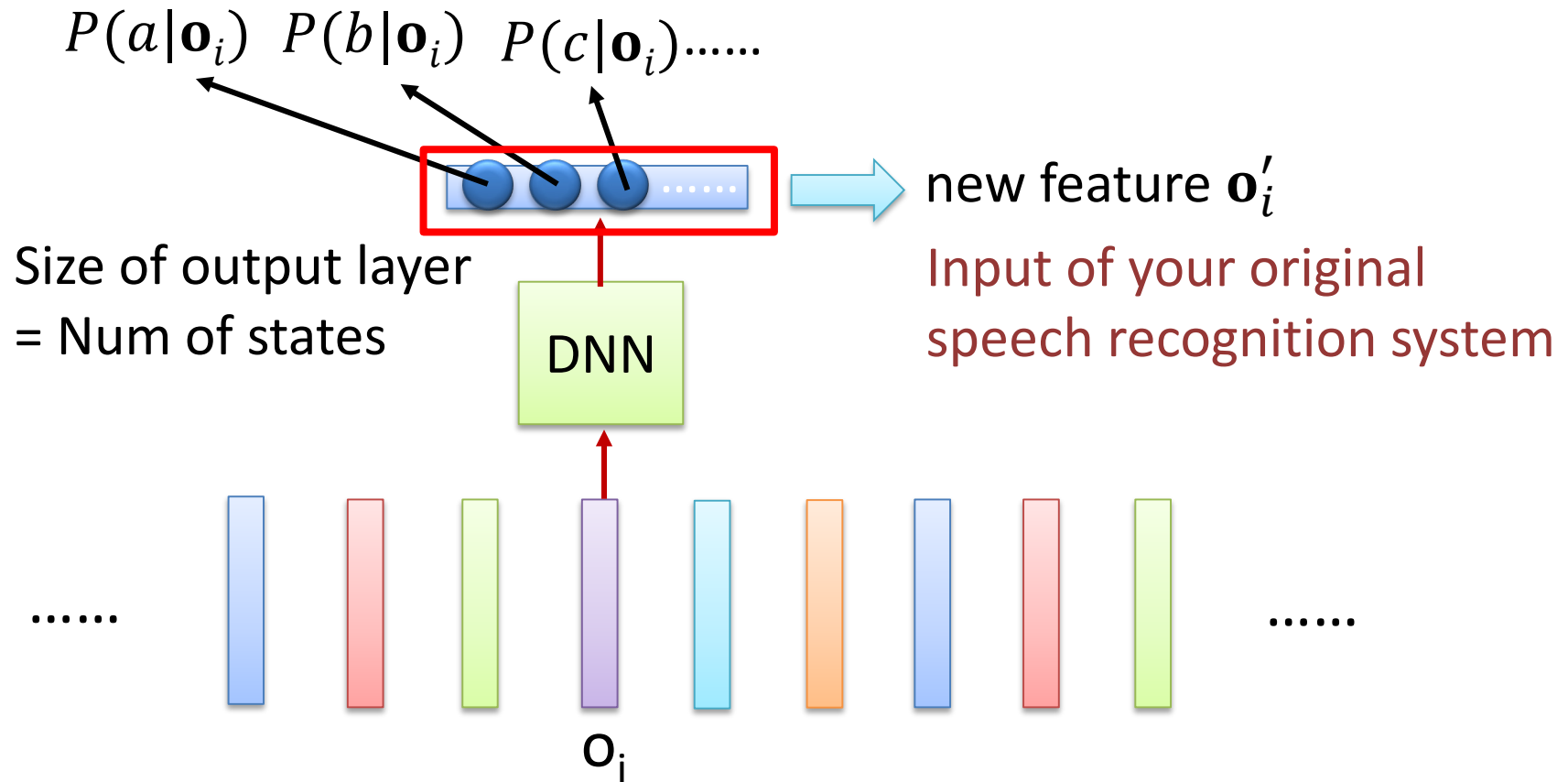


Efforts for exploiting deep learning

How to use Deep Learning?

Way 1: Tandem





Last hidden layer or bottleneck layer are also possible.

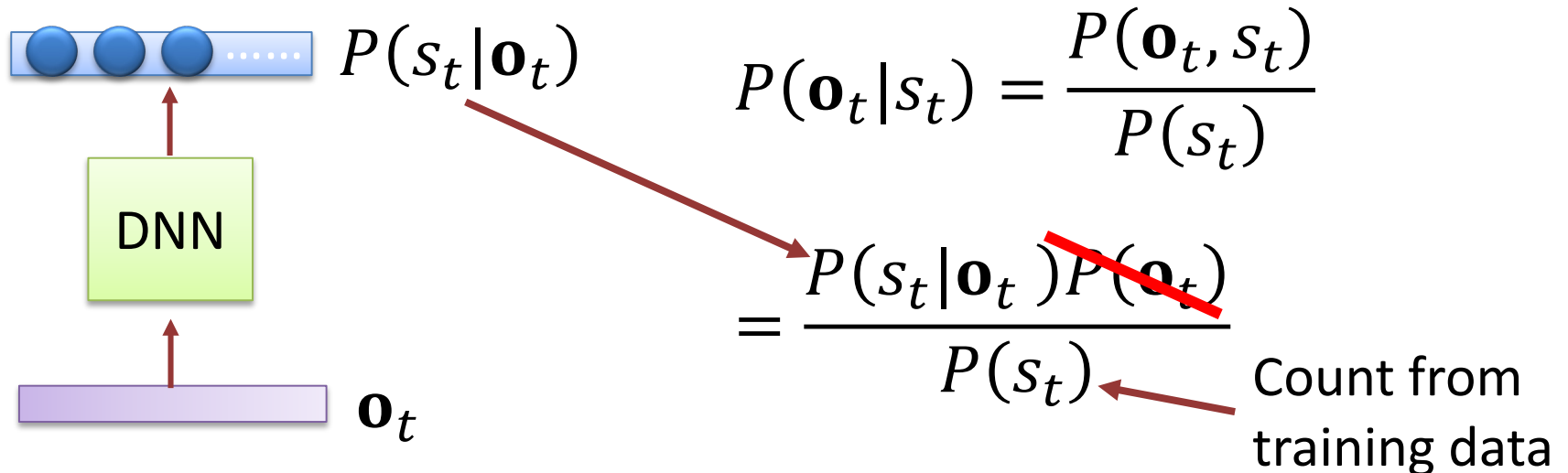
How to use Deep Learning?

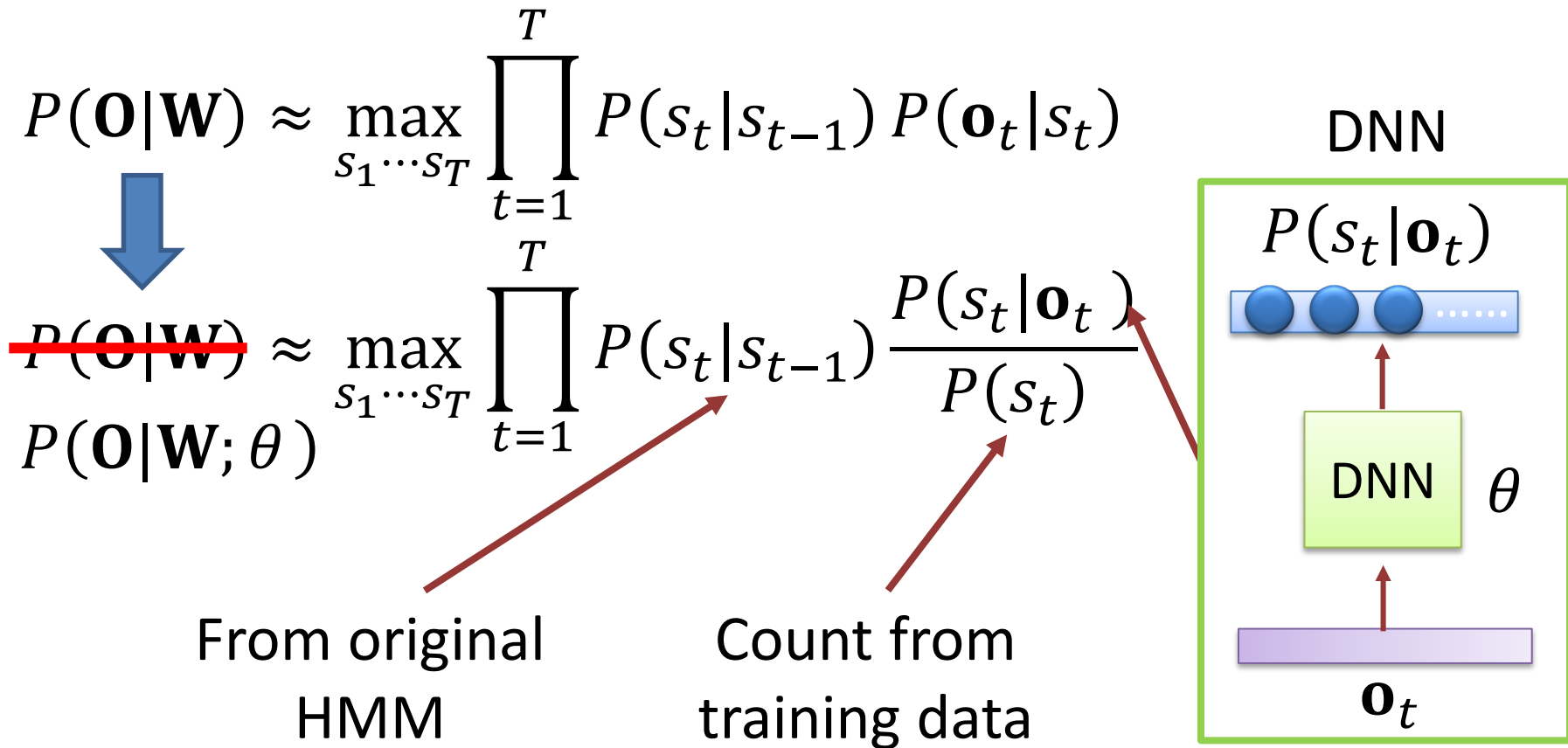
Way 2: DNN-HMM hybrid



$$\mathbf{W}^* = \arg \max_{\mathbf{W} \in L} p(\mathbf{W}|\mathbf{O}) = \arg \max_{\mathbf{W} \in L} p(\mathbf{O}|\mathbf{W})P(\mathbf{W})$$

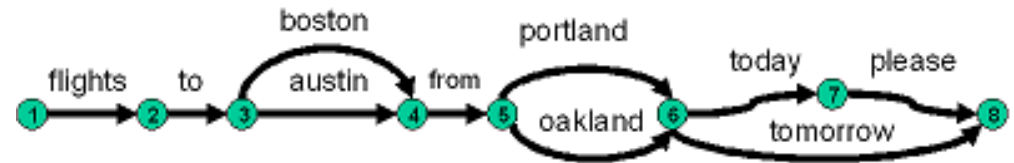
$$P(\mathbf{O}|\mathbf{W}) \approx \max_{s_1 \dots s_T} \prod_{t=1}^T P(s_t|s_{t-1}) \underbrace{P(\mathbf{o}_t|s_t)}_{\text{From DNN}}$$





This assembled vehicle works

◆ Sequential Training



$$\mathbf{W}^* = \arg \max_{\mathbf{W} \in L} p(\mathbf{O} | \mathbf{W}; \theta) P(\mathbf{W})$$

Given training data $(\mathbf{O}_1, \mathbf{W}_1^*), (\mathbf{O}_2, \mathbf{W}_2^*), \dots (\mathbf{O}_r, \mathbf{W}_r^*), \dots$

Find-tune the DNN parameters θ such that

$$P(\mathbf{O}_r | \mathbf{W}_r^*; \theta) P(\mathbf{W}_r^*) \rightarrow \text{increase}$$

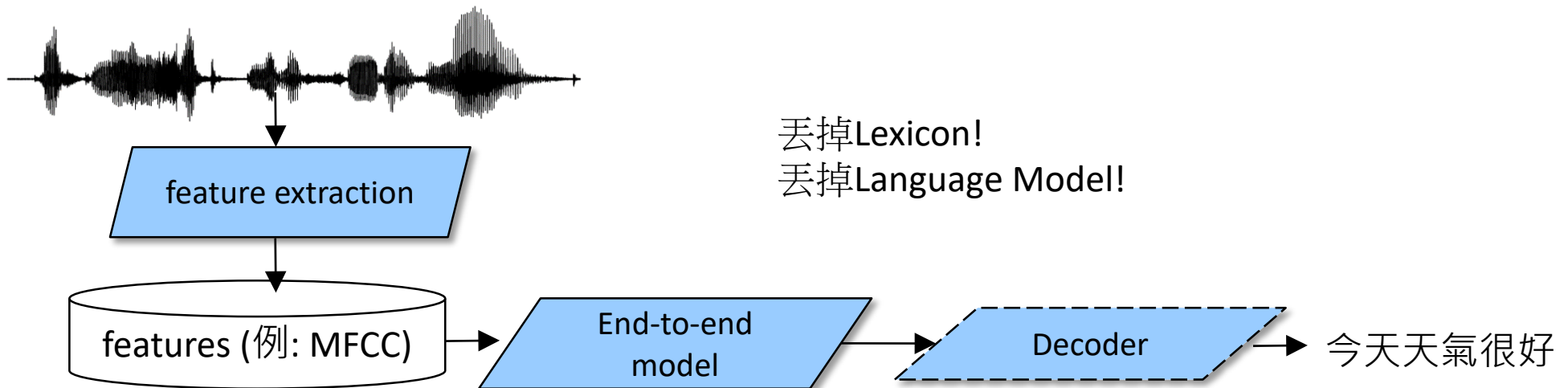
$$P(\mathbf{O}_r | \mathbf{W}; \theta) P(\mathbf{W}) \rightarrow \text{decrease}$$

(\mathbf{W} is any word sequence different from \mathbf{W}_r^*)

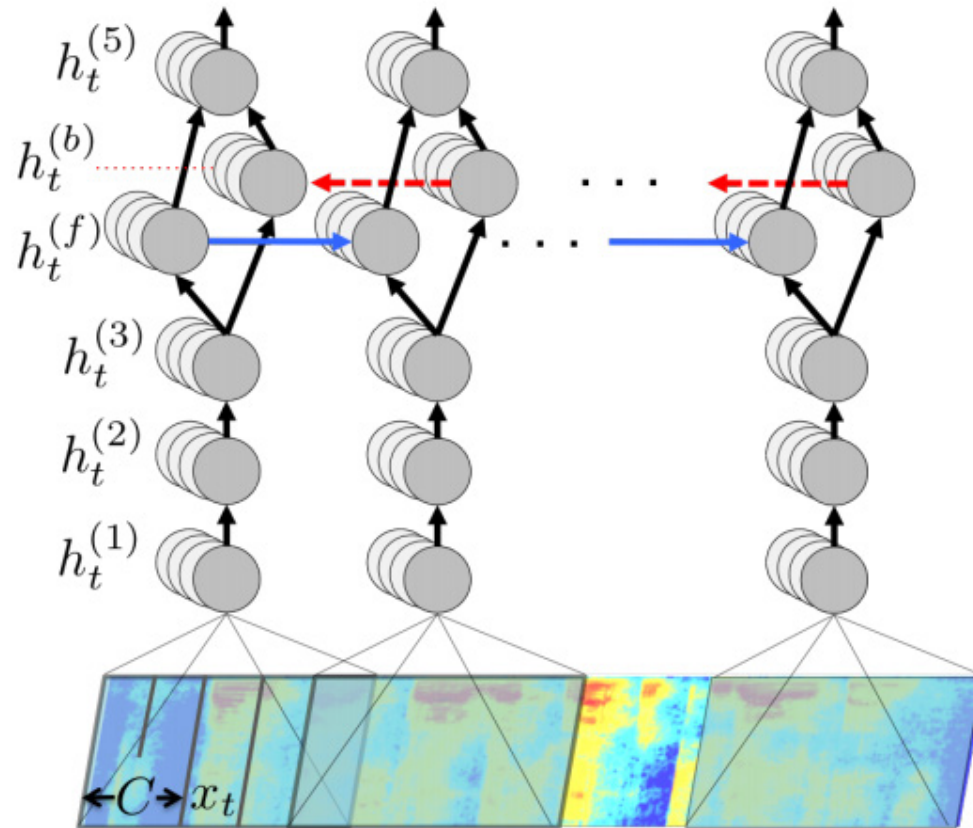
How to use Deep Learning?

Way 3: End-to-end





- ◆ **Directly train model to solve task (“speech-to-text”)**
 - single model trained
 - no separate acoustic and language models
- ◆ **More complicated to incorporate additional LM data**
- ◆ **Output layer從state換成character, phone or words**



Input: acoustic features
(spectrograms)

Output: characters
(and space)
+ null (\sim)

No phoneme and lexicon
(No OOV problem)

HIS FRIEND'S

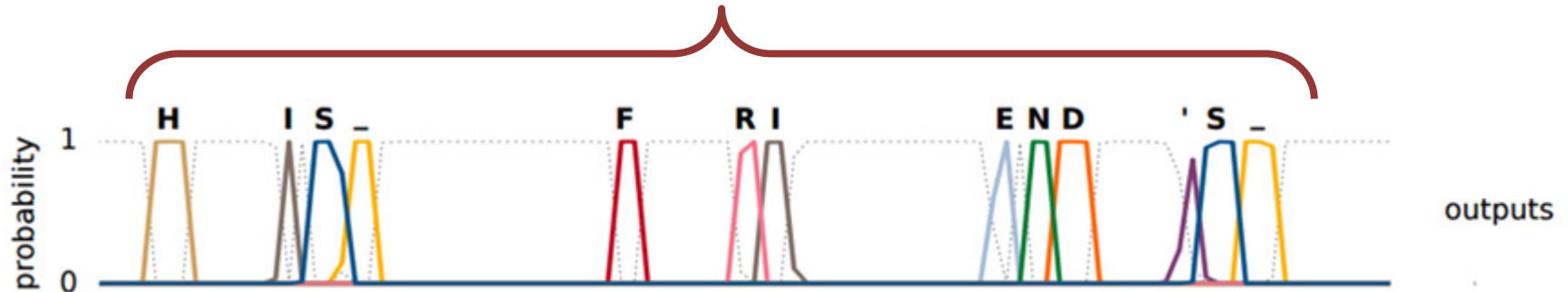


Figure 4. Network outputs. The figure shows the frame-level character probabilities emitted by the CTC layer (different colour for each character, dotted grey line for 'blanks'), along with the corresponding training errors, while processing an utterance. The target transcription was 'HIS_FRIENDS_', where the underscores are end-of-word markers. The network was trained with WER loss, which tends to give very sharp output decisions, and hence sparse error signals (if an output probability is 1, nothing else can be sampled, so the gradient is 0 even if the output is wrong). In this case the only gradient comes from the extraneous apostrophe before the 'S'. Note that the characters in common sequences such as 'IS', 'RI' and 'END' are emitted very close together, suggesting that the network learns them as single sounds.

Graves, Alex, and Navdeep Jaitly. "Towards end-to-end speech recognition with recurrent neural networks." *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*. 2014.

- ◆ training corpus: 1.2 billions words, vocabulary: 1.7 million words.
- ◆ 125,000 hours of semi-supervised acoustic training data
 - deep bi-directional LSTM RNNs with CTC loss

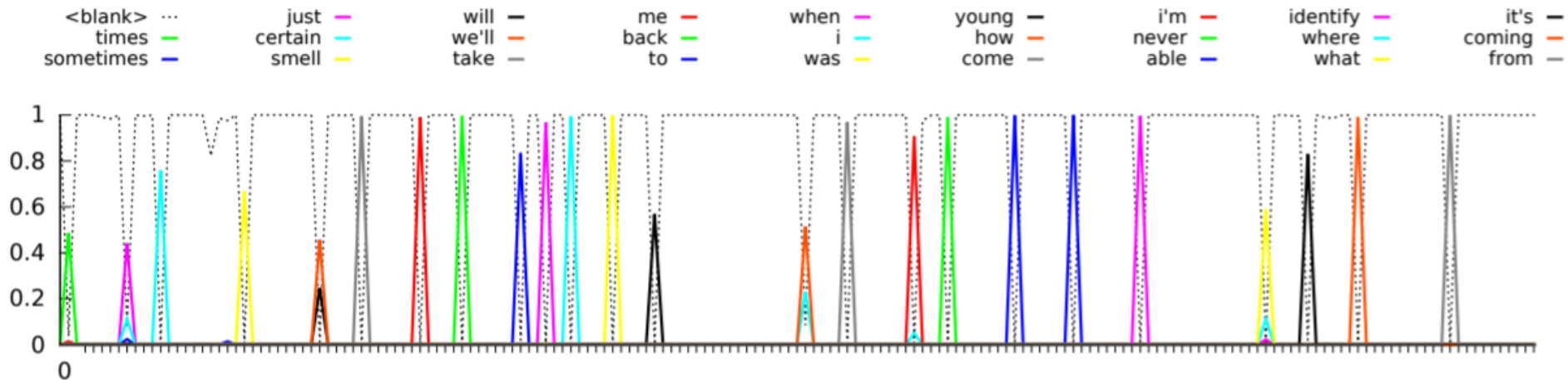


Figure 1: The word posterior probabilities as predicted by the NSR model at each time-frame (30 msec) for a segment of music video ‘Stressed Out’ by Twenty One Pilots. We only plot the word with highest posterior and the missing words from the correct transcription: ‘*Sometimes a certain smell will take me back to when I was young, how come I’m never able to identify where it’s coming from*’.

Why Deep Learning?



Layer X Size	Word Error Rate (%)
1 X 2k	24.2
2 X 2k	20.4
3 X 2k	18.4
4 X 2k	17.8
5 X 2k	17.2
7 X 2k	17.1

Not surprised, more parameters, better performance

◆ Example Architecture from Google (2015)

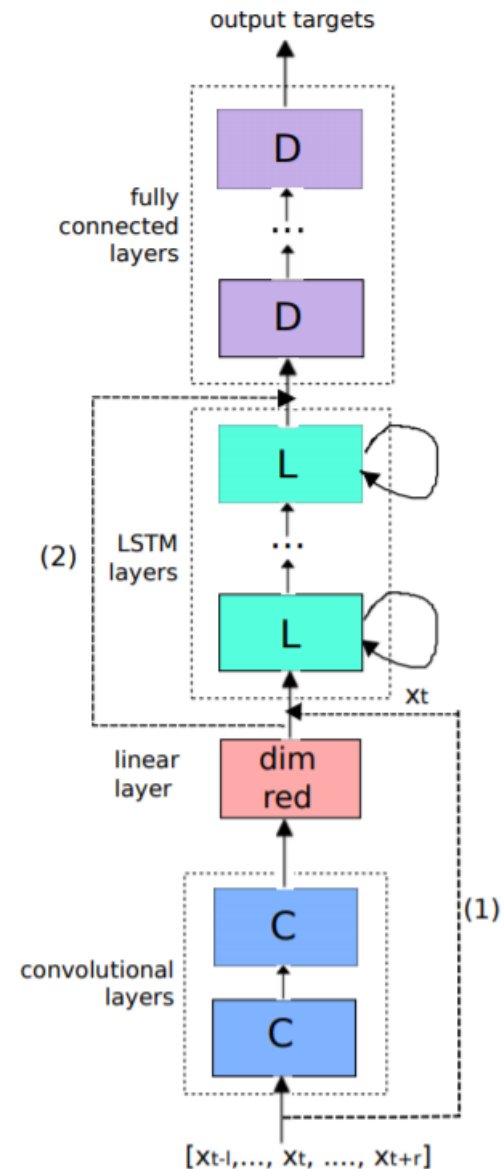
- C: CNN layer (with pooling)
- L: LSTM layer
- D: fully connected layer

◆ Two multiple layer “skips”

- (1) connects input to LSTM input
- (2) connects CNN output to DNN input

◆ Additional linear projection layer

- reduces dimensionality
- and number of network parameters!

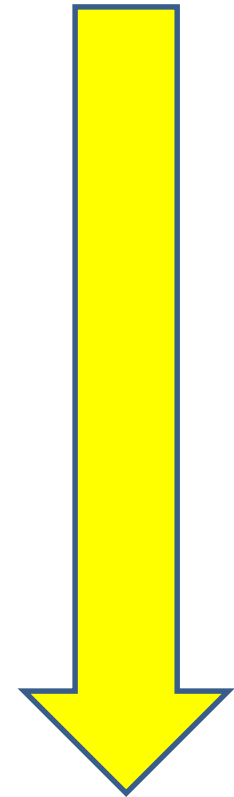
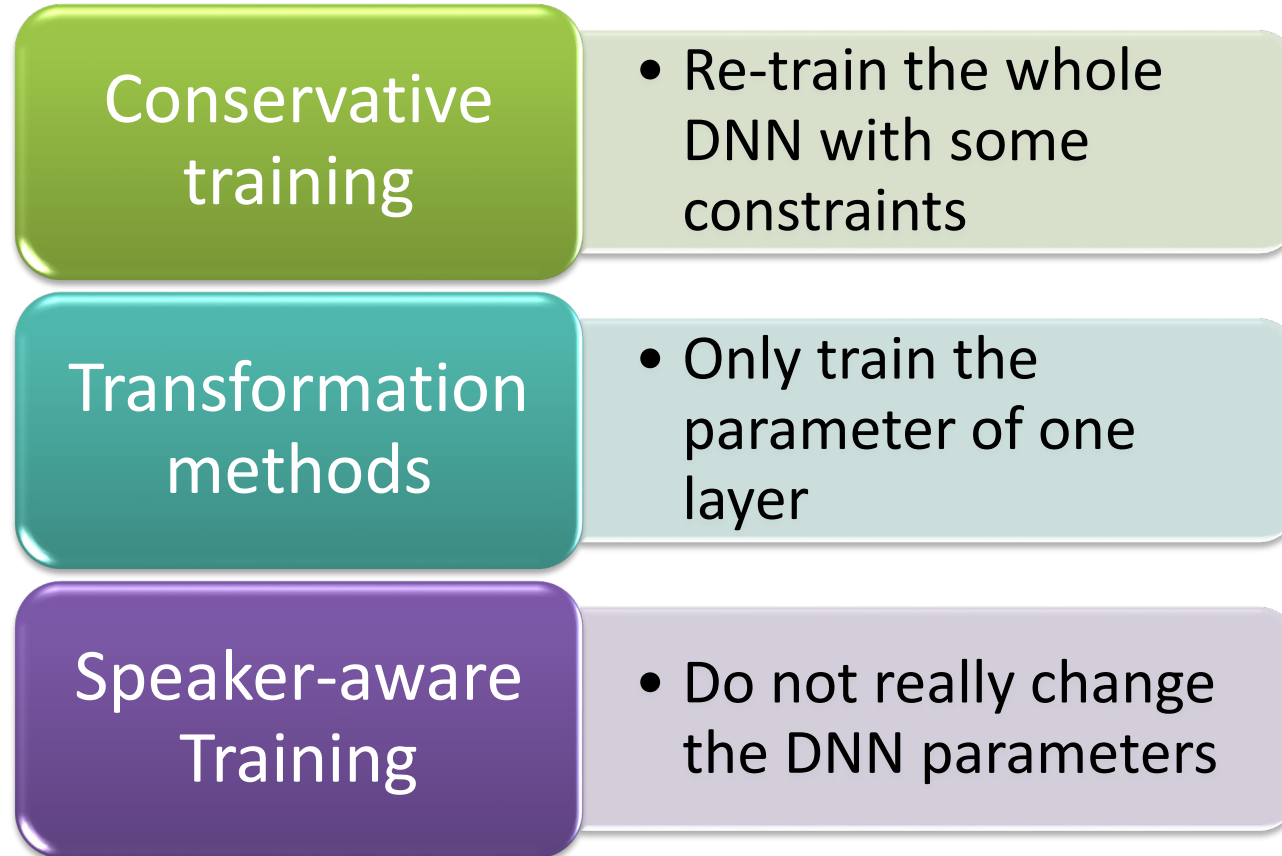


Speaker Adaptation

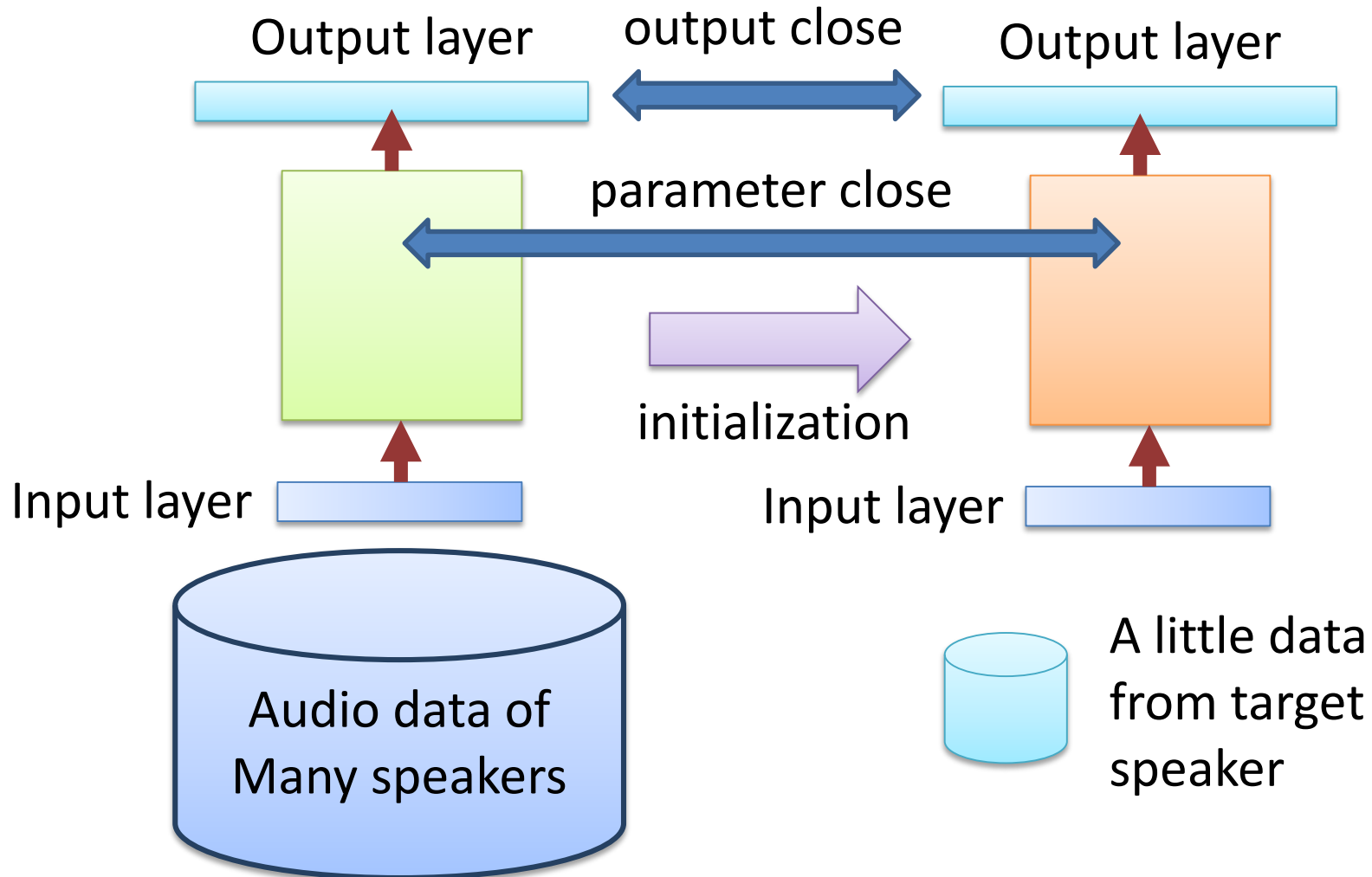


- ◆ **Speaker adaptation: use different models to recognition the speech of different speakers**
 - **Collect the audio data of each speaker**

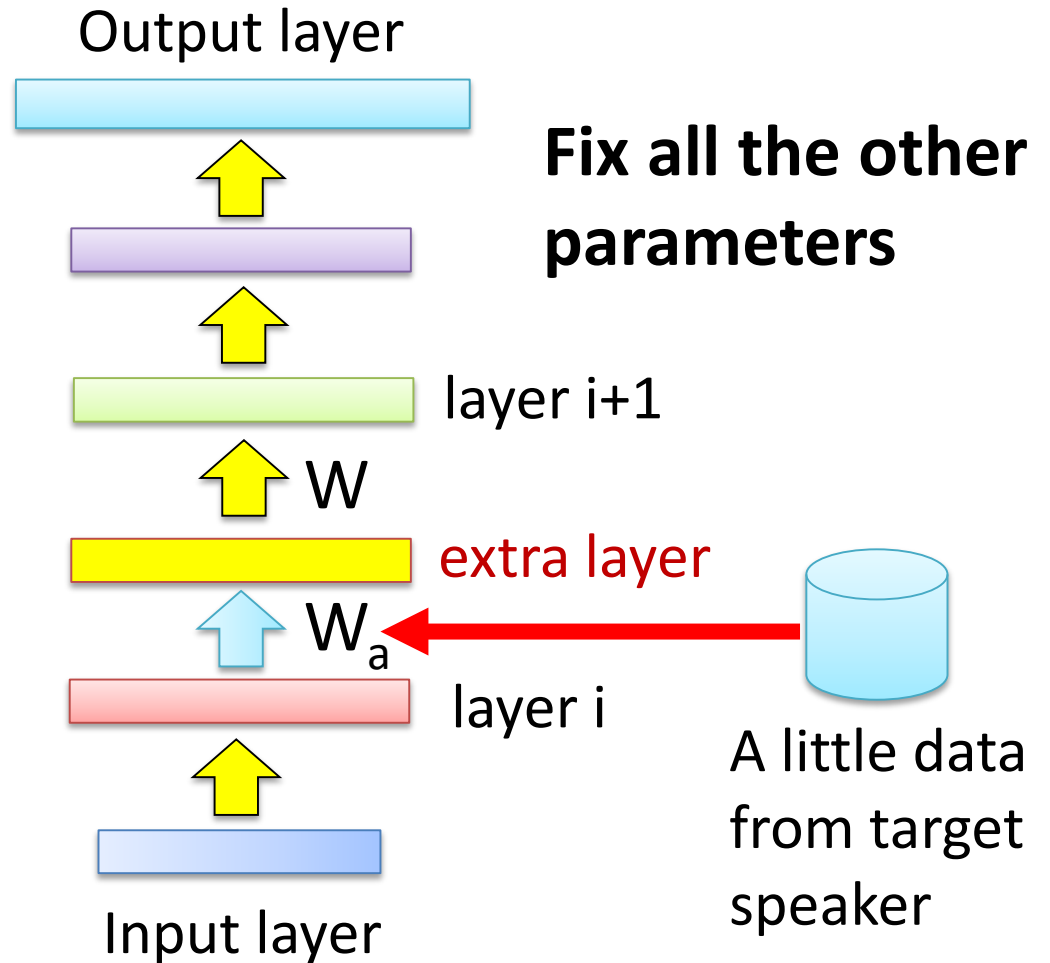
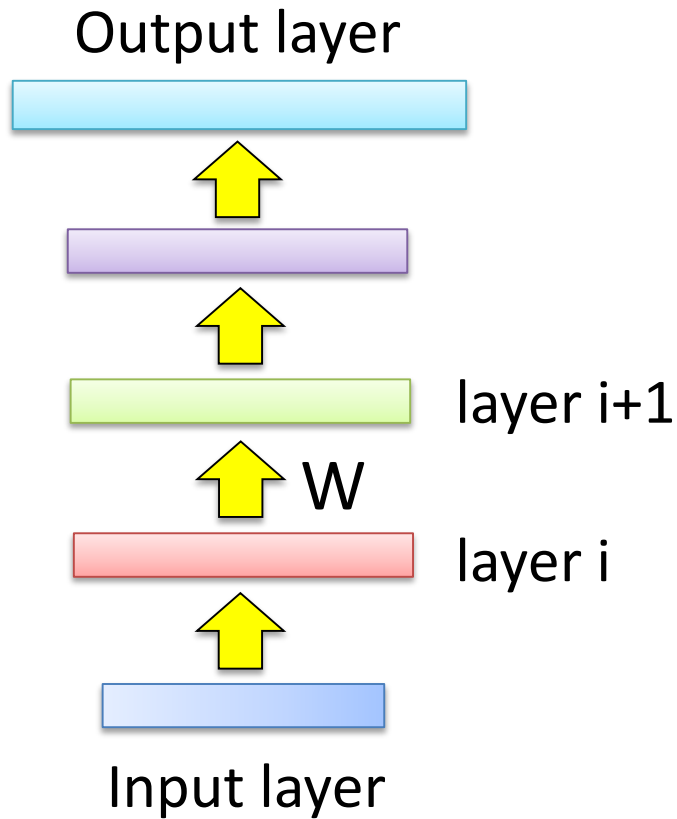
- ◆ **A DNN model for each speaker**
 - **Challenge: limited data for training**
 - * Not enough data for directly training a DNN model
 - * Not enough data for just fine-tune a speaker independent DNN model



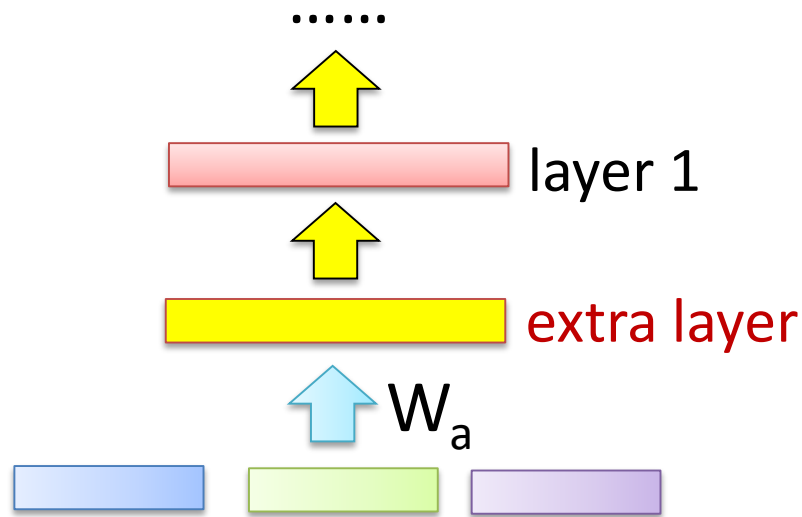
Need less training data



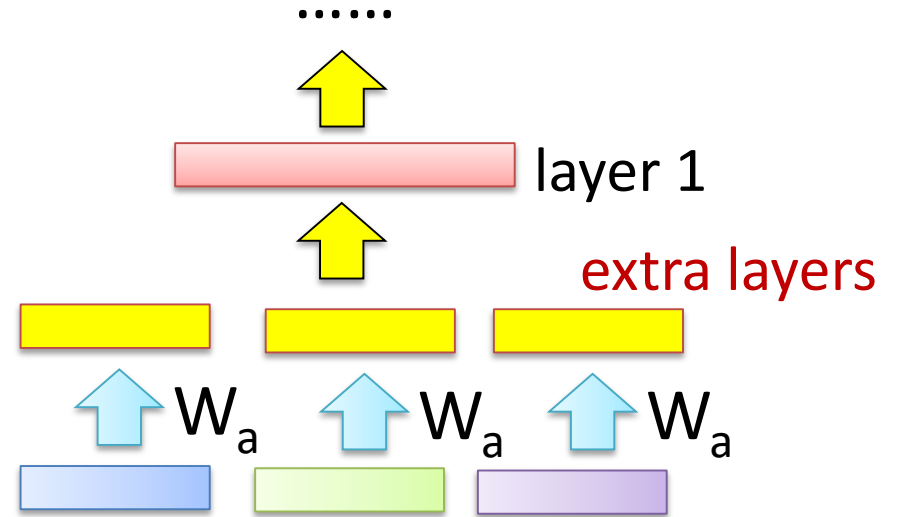
Add an extra layer



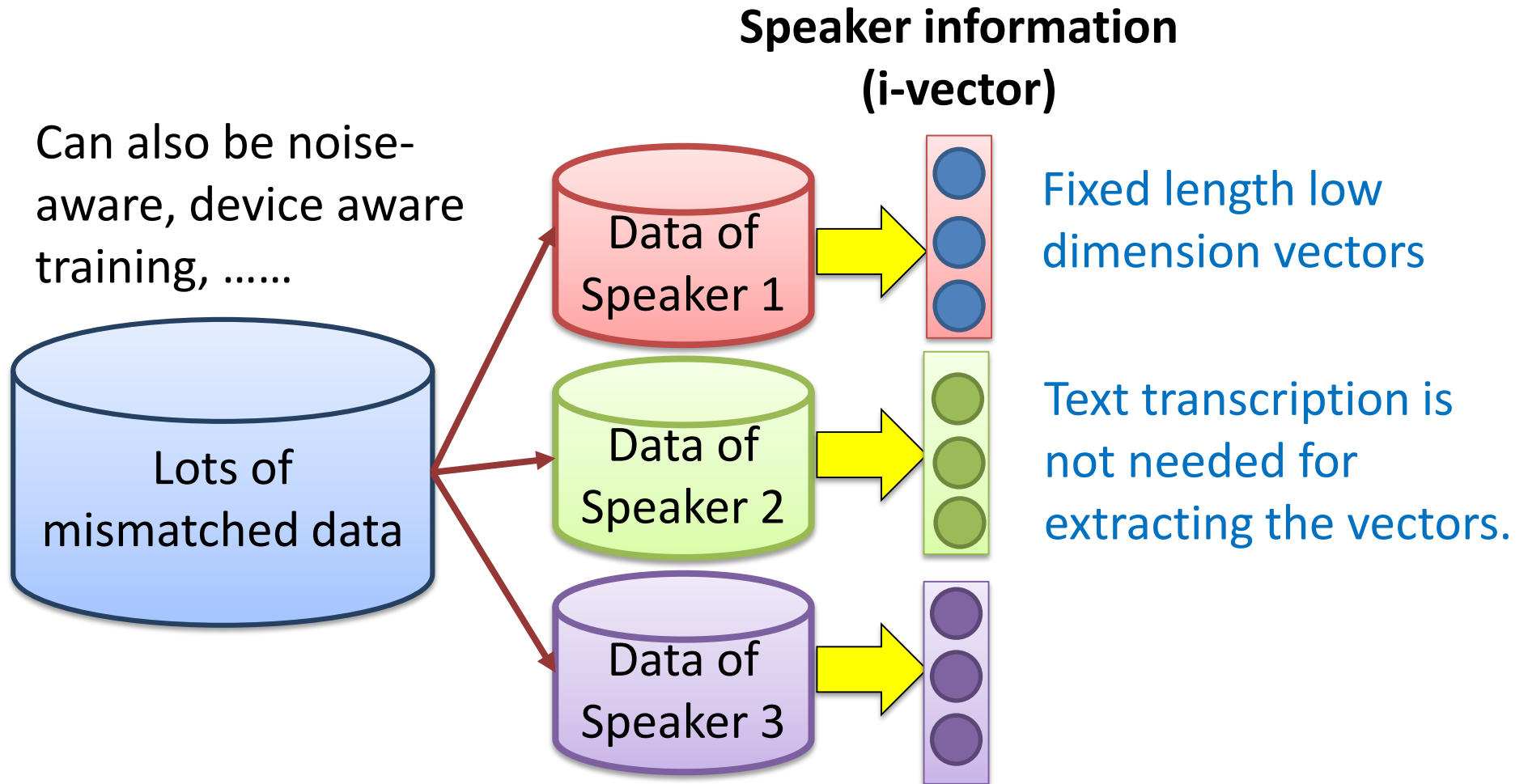
- ◆ Add the extra layer between the input and first layer
- ◆ With splicing



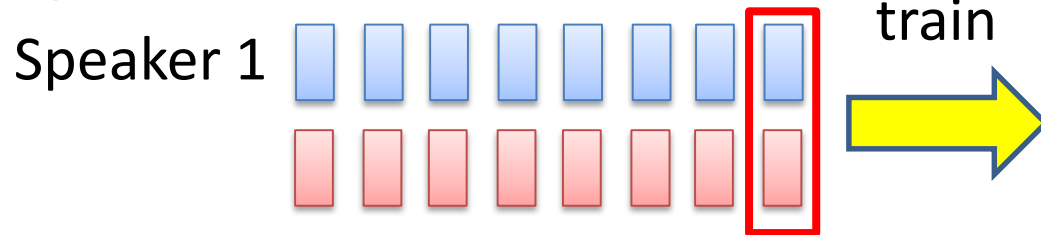
Larger W_a → More data



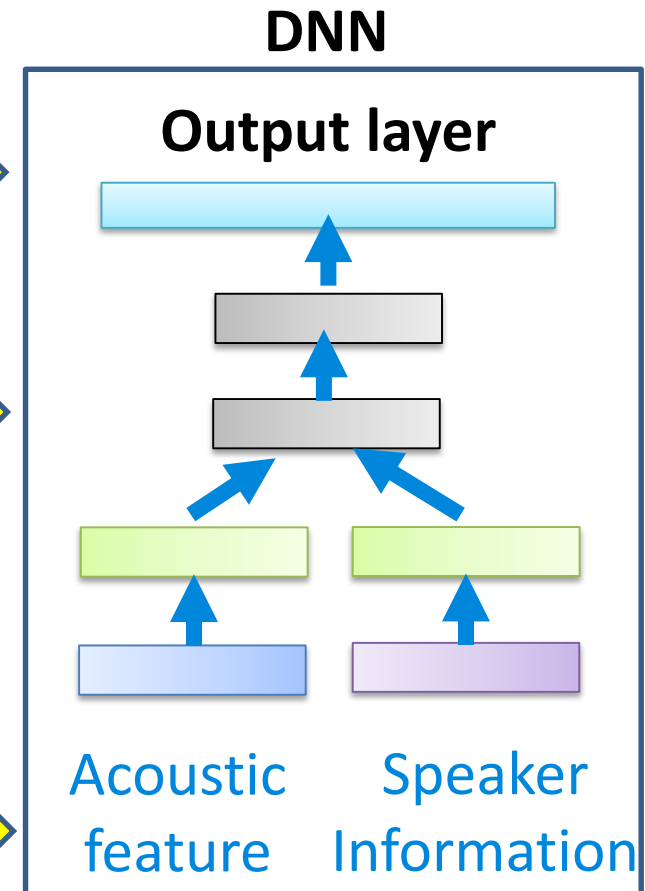
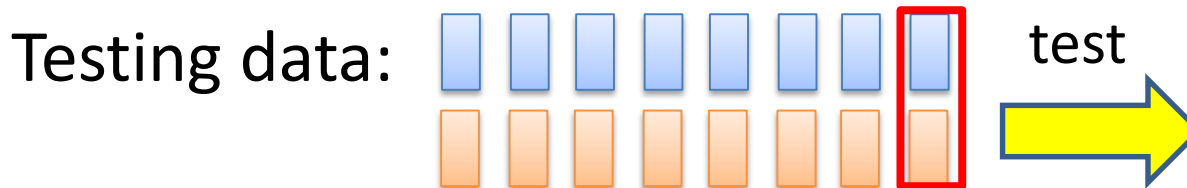
Smaller W_a → less data



Training data:



Acoustic features are appended with speaker information features



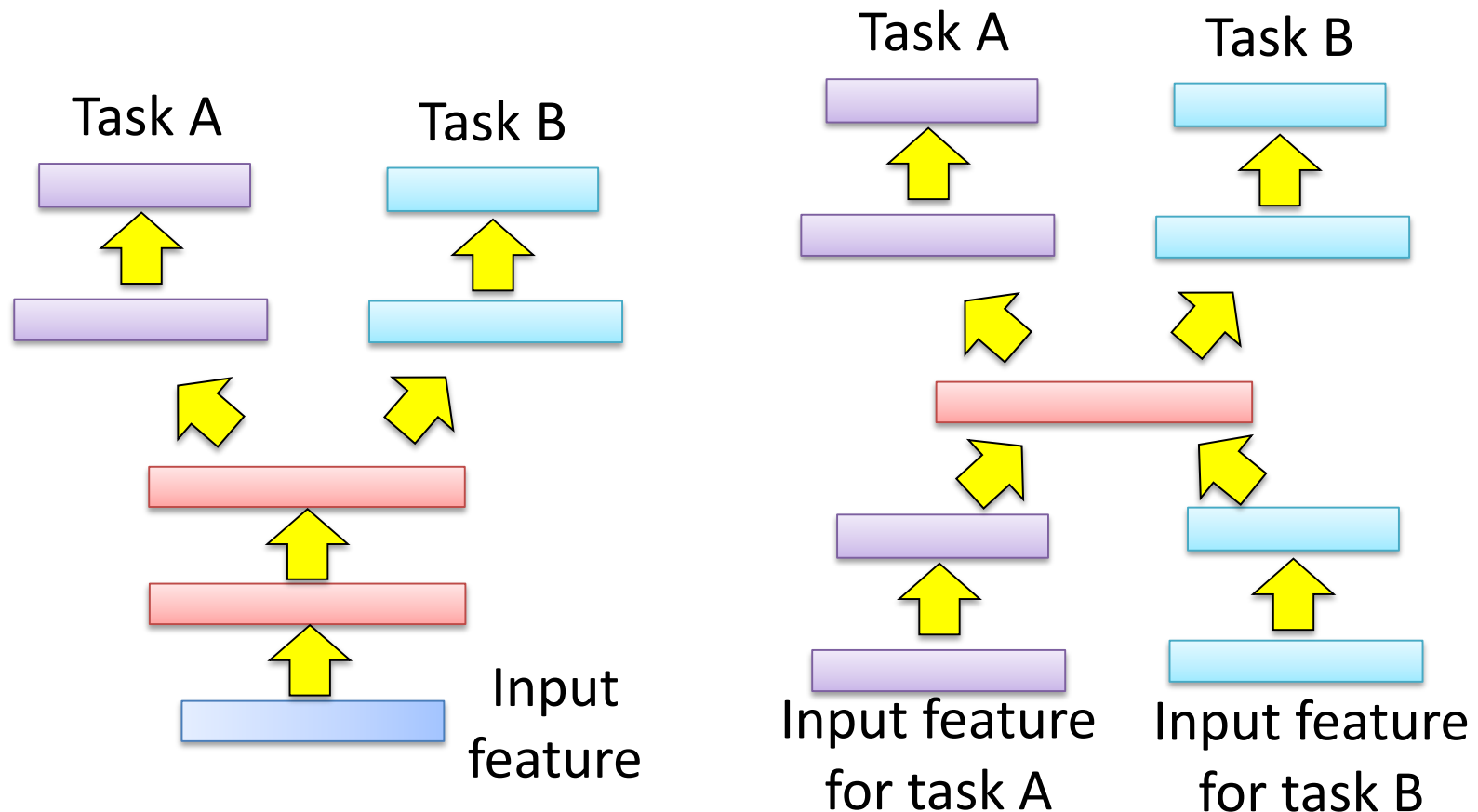
All the speaker use the same DNN model

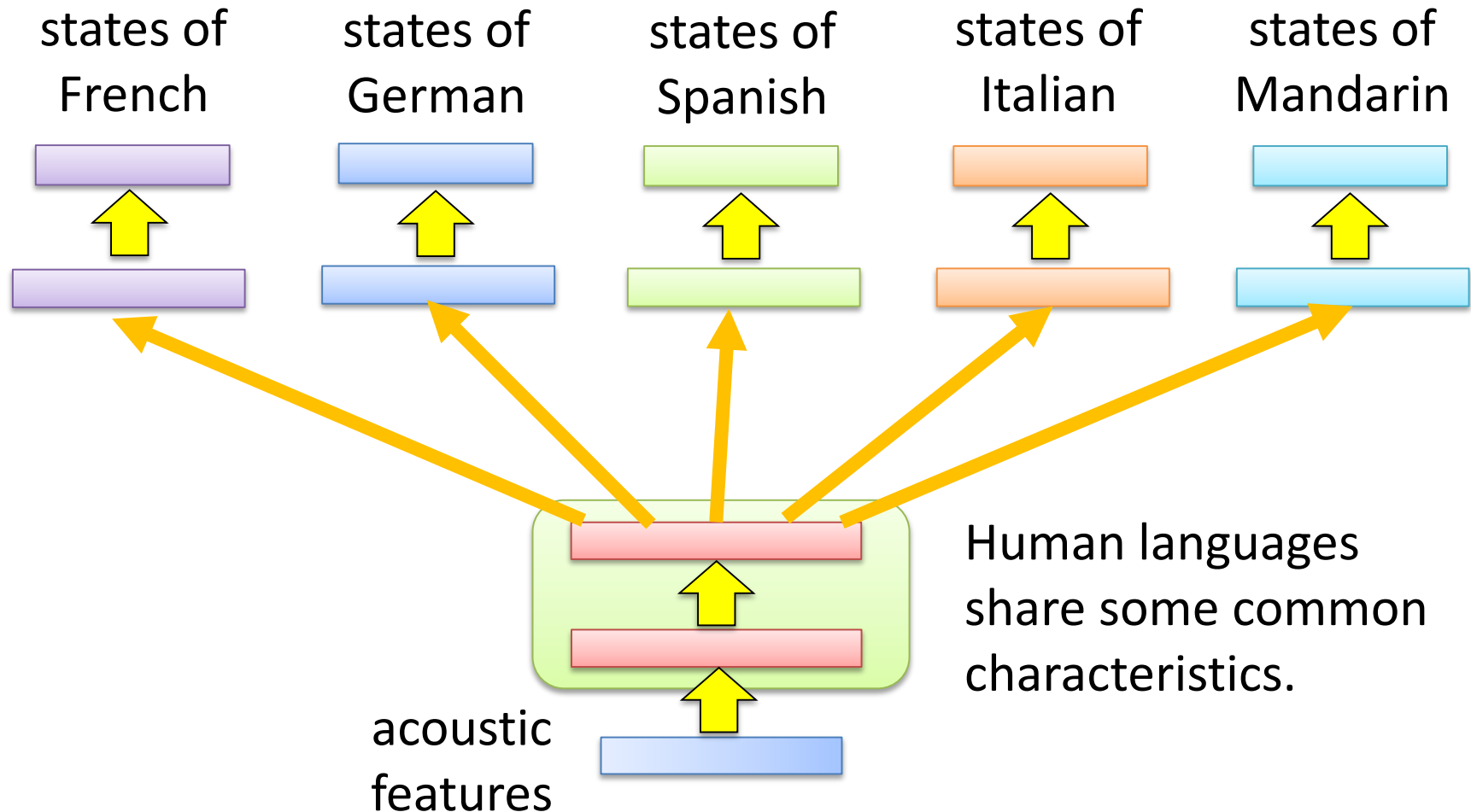
Different speaker augmented by different features

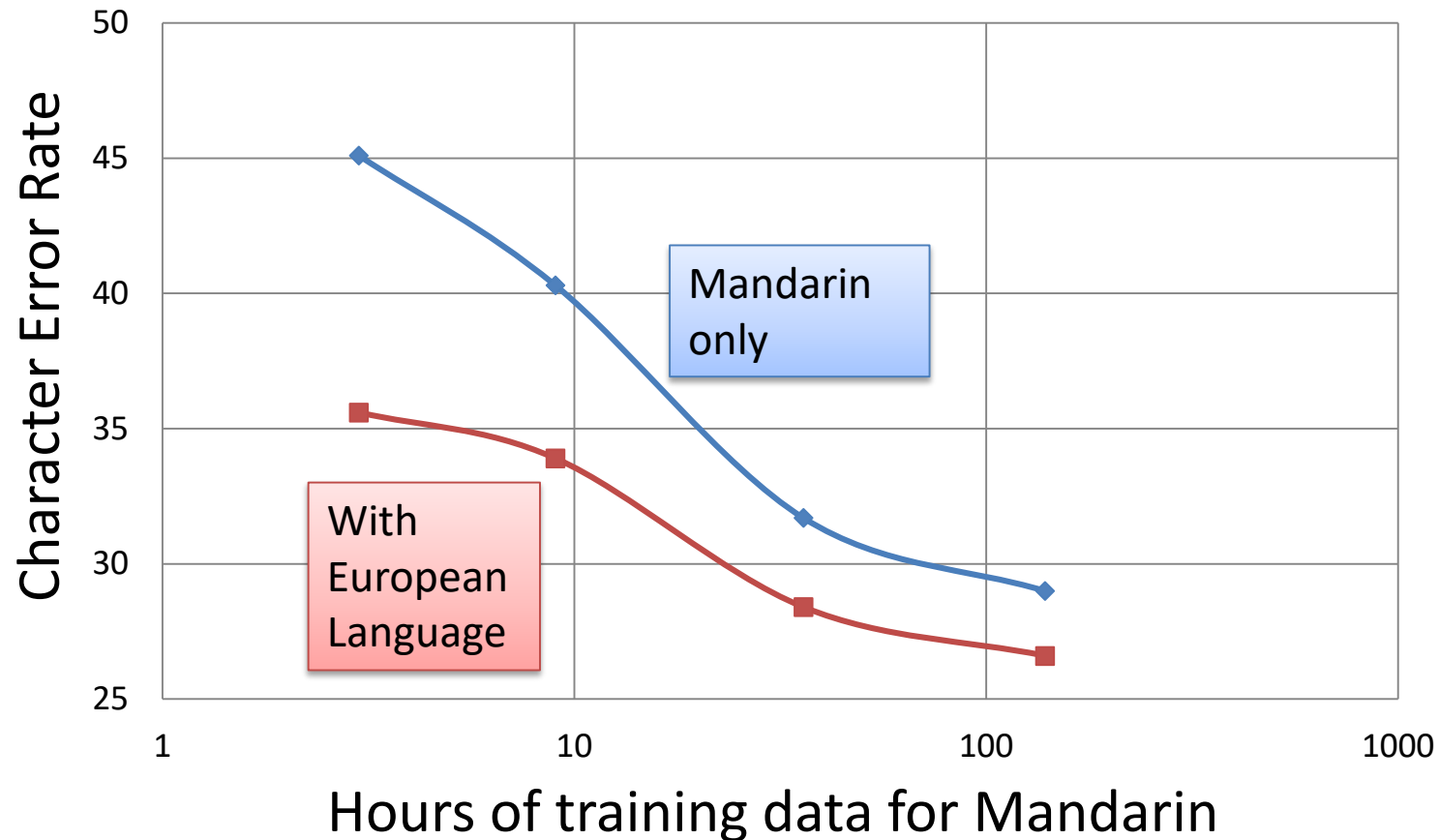
Multi-task Learning

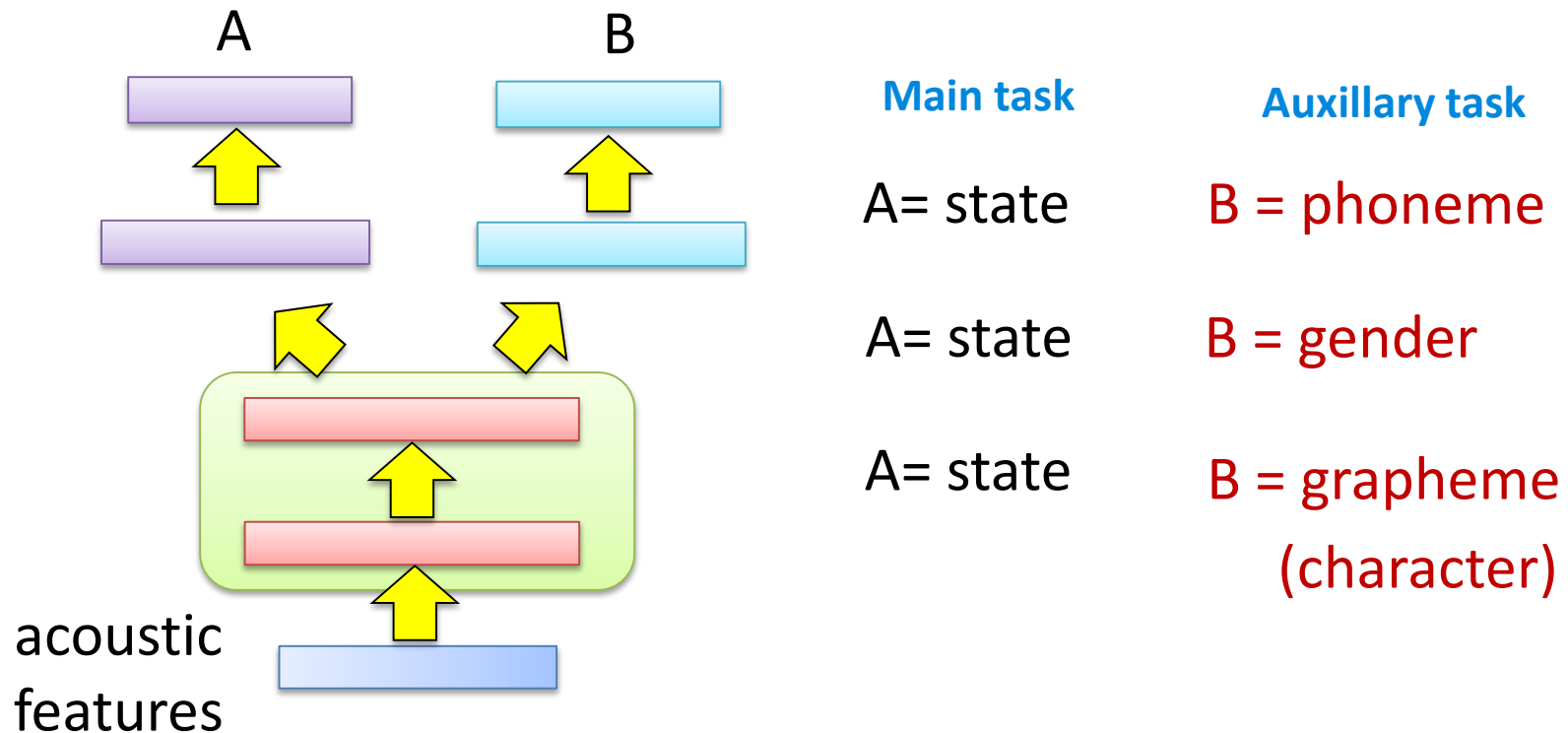


- ◆ The multi-layer structure makes DNN suitable for multitask learning





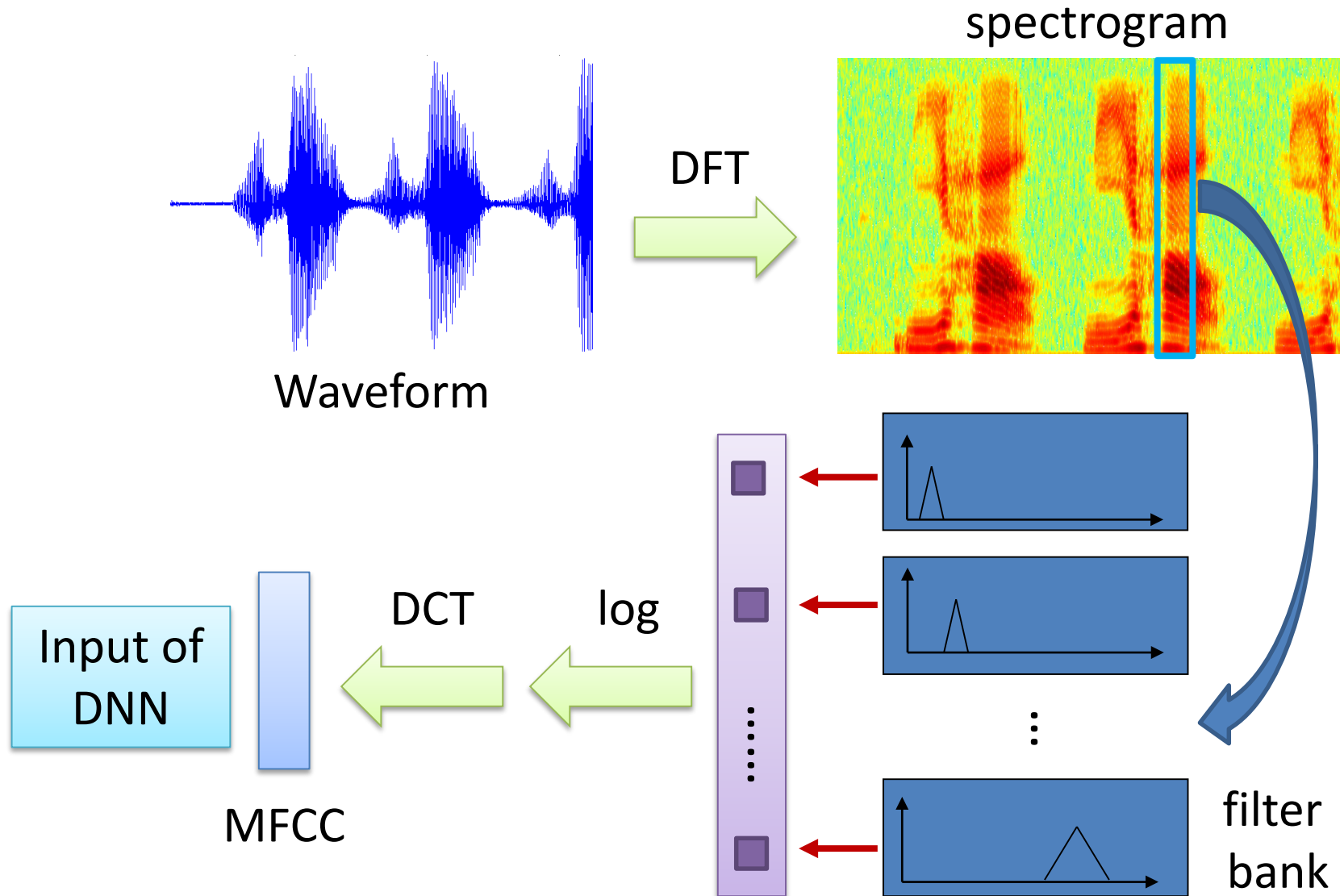


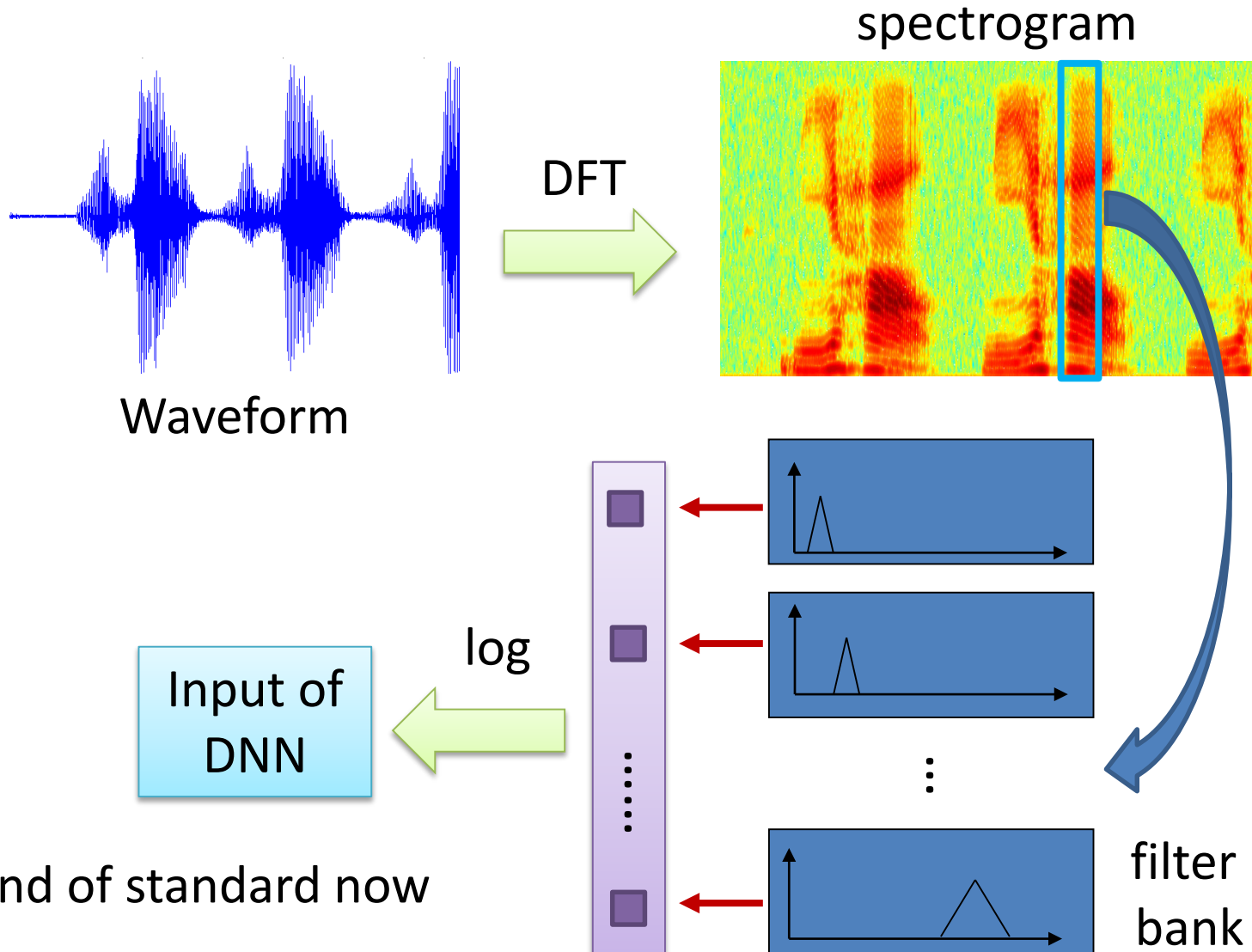


Deep Learning for Acoustic Modeling

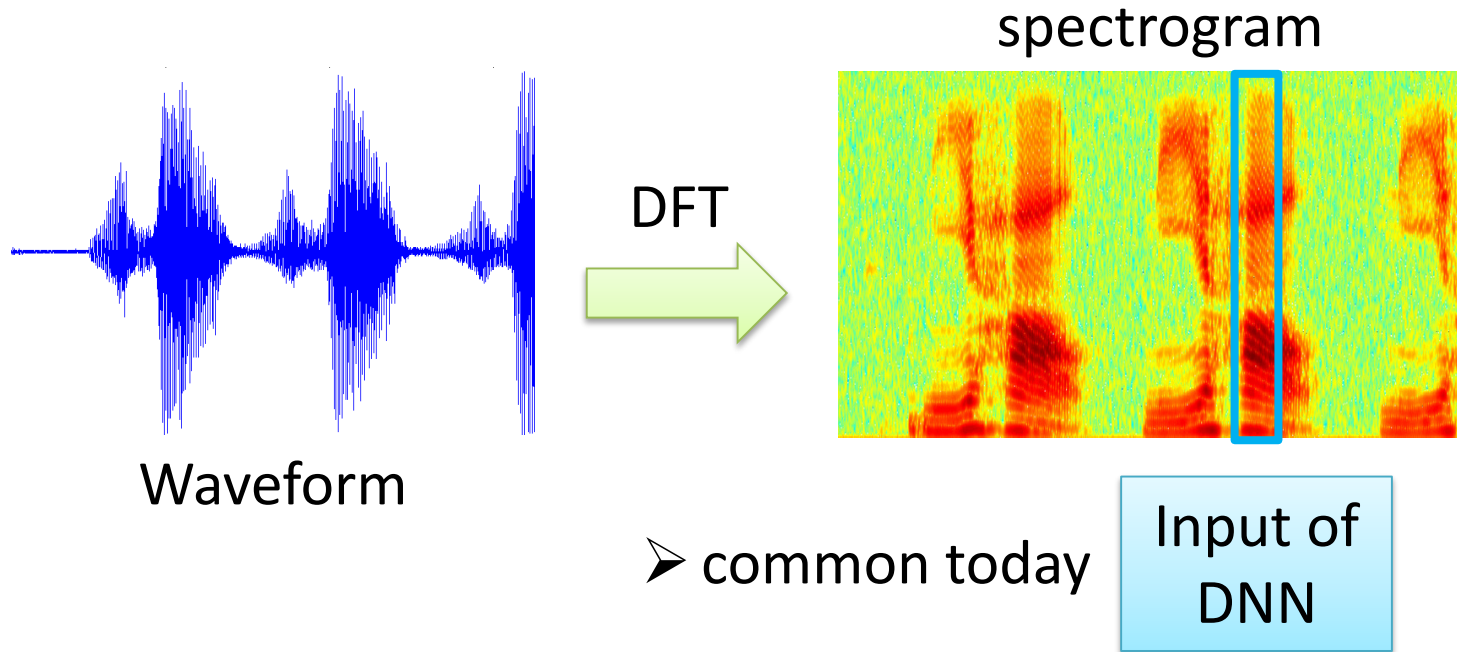
New acoustic features





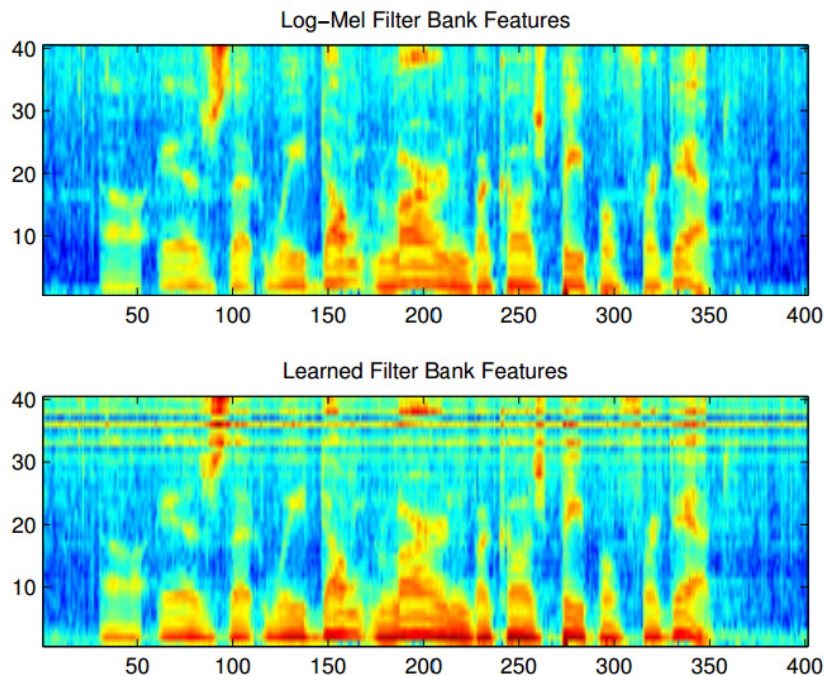
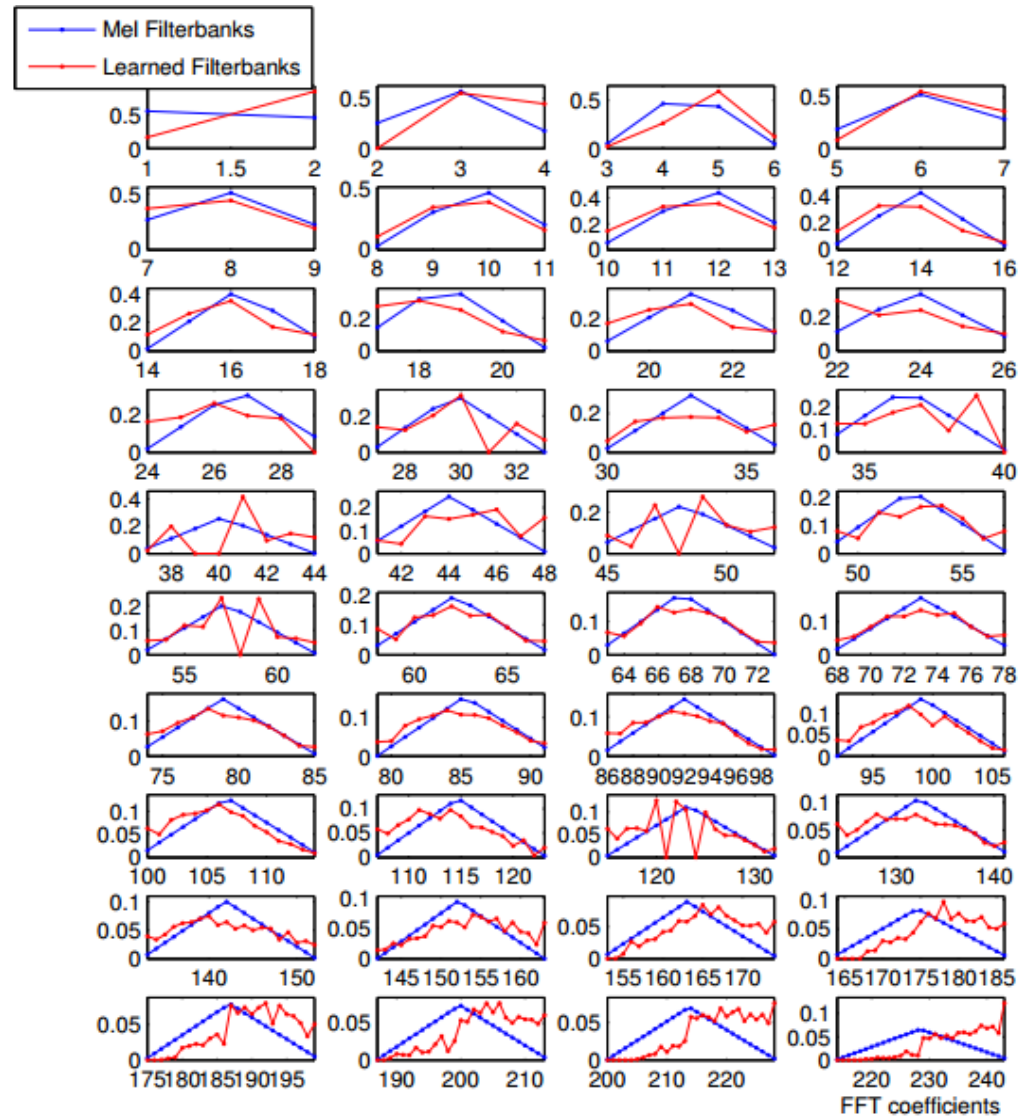


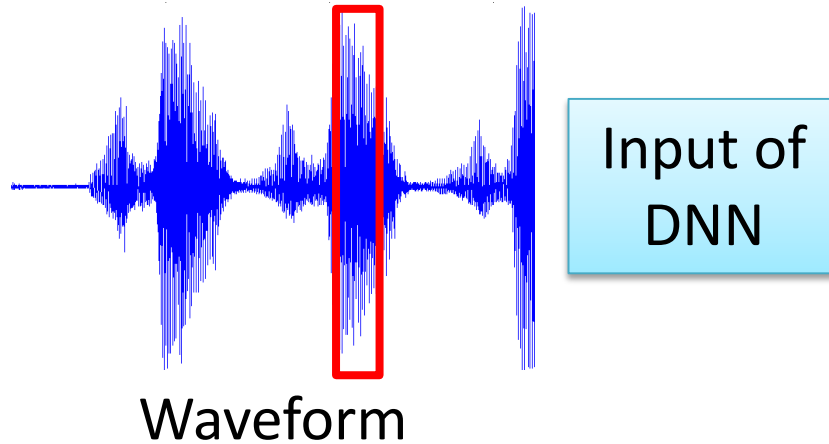
➤ Kind of standard now




- 5% relative improvement over filterbank output

◆ Learning fbanks within DNN



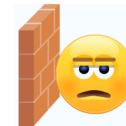


➤ If success, no Signal & Systems 

➤ People tried, but **not** better than spectrogram yet

Tüske, Z et al., "Acoustic modeling with deep neural networks using raw time signal for LVCSR," In *INTERPSEECH 2014*

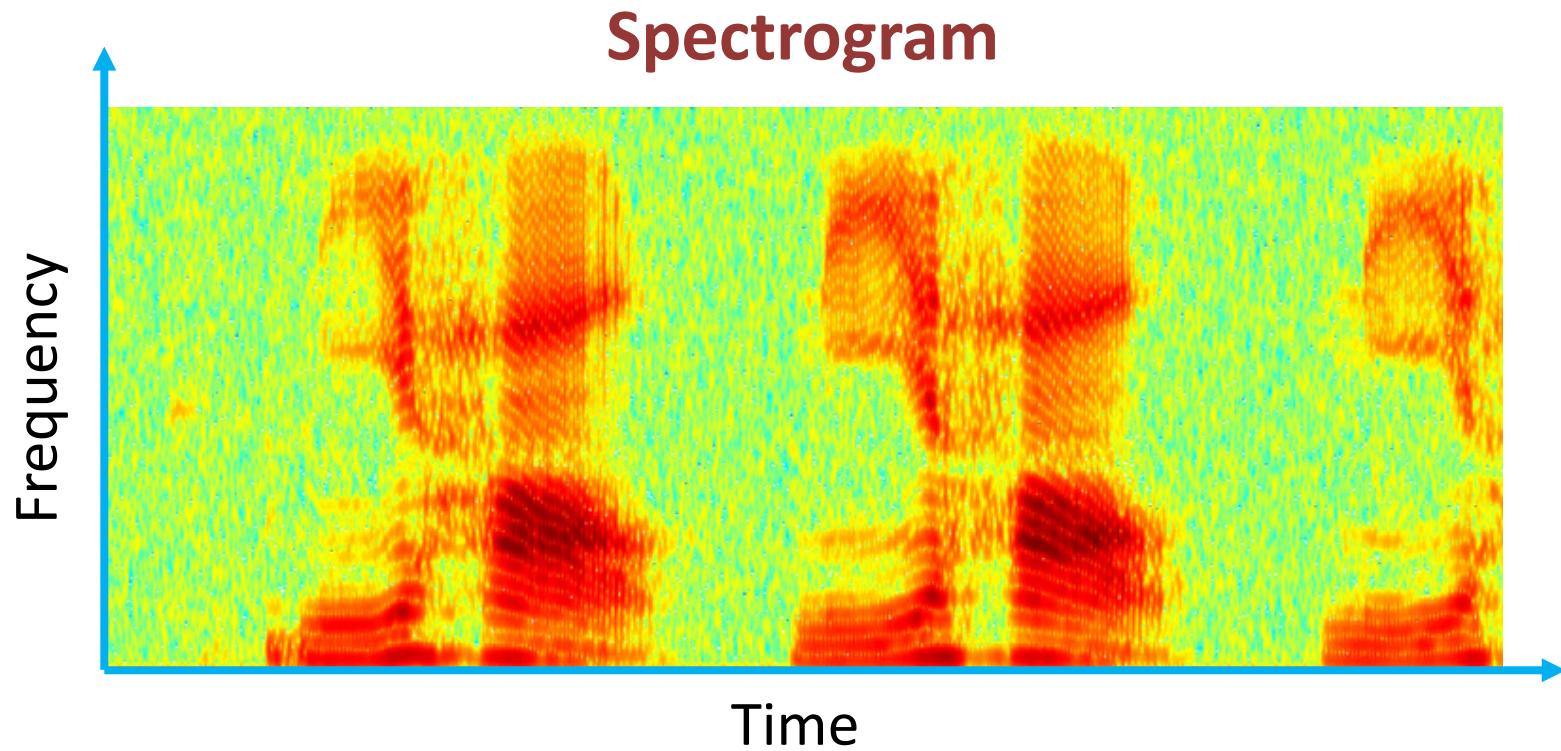
➤ Still need to take Signal & Systems

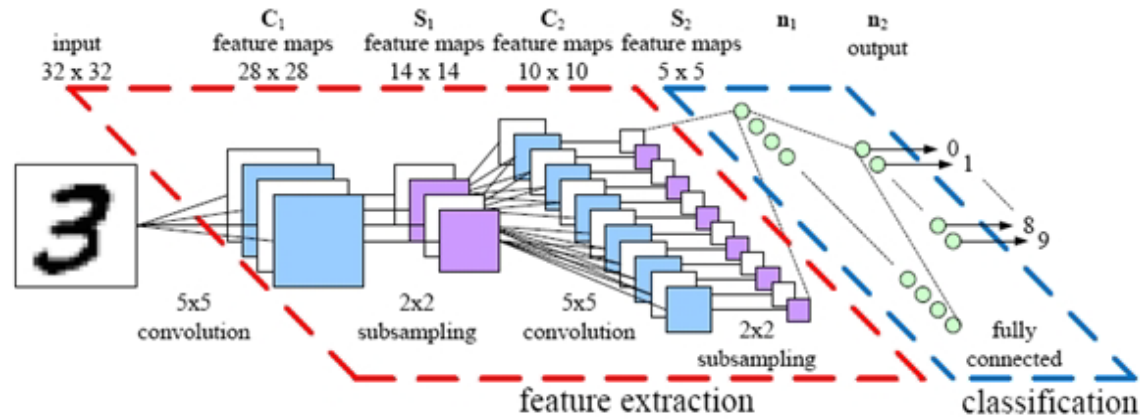


Convolutional Neural Network (CNN)

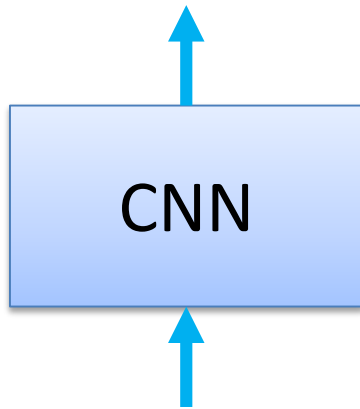


- ◆ Speech can be treated as images

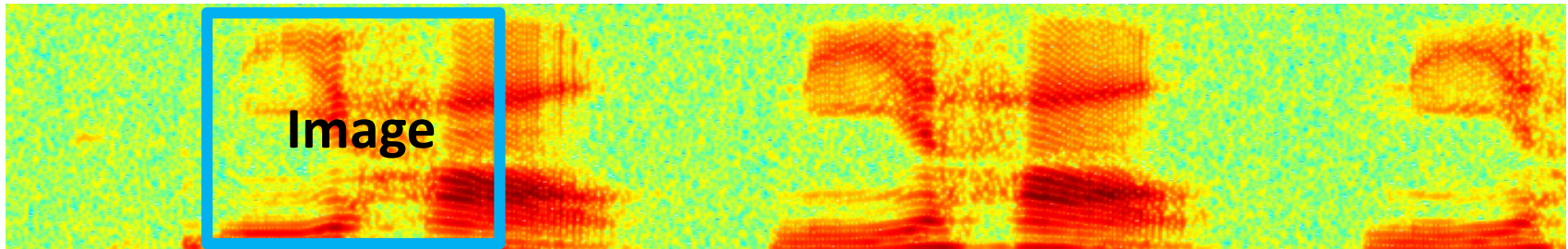


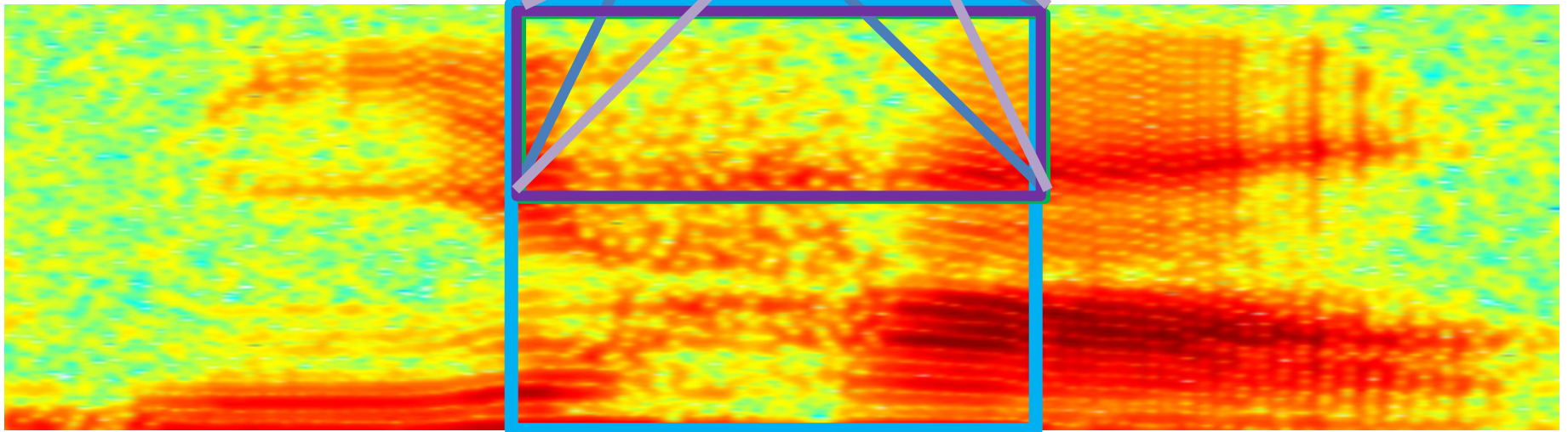
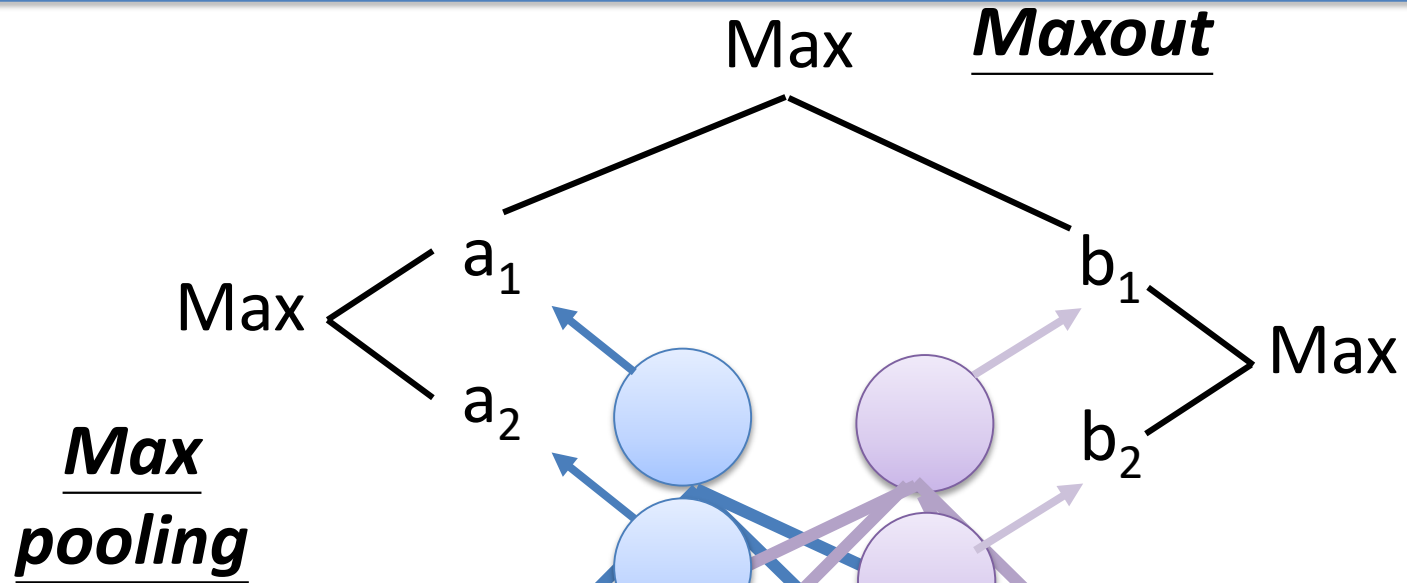


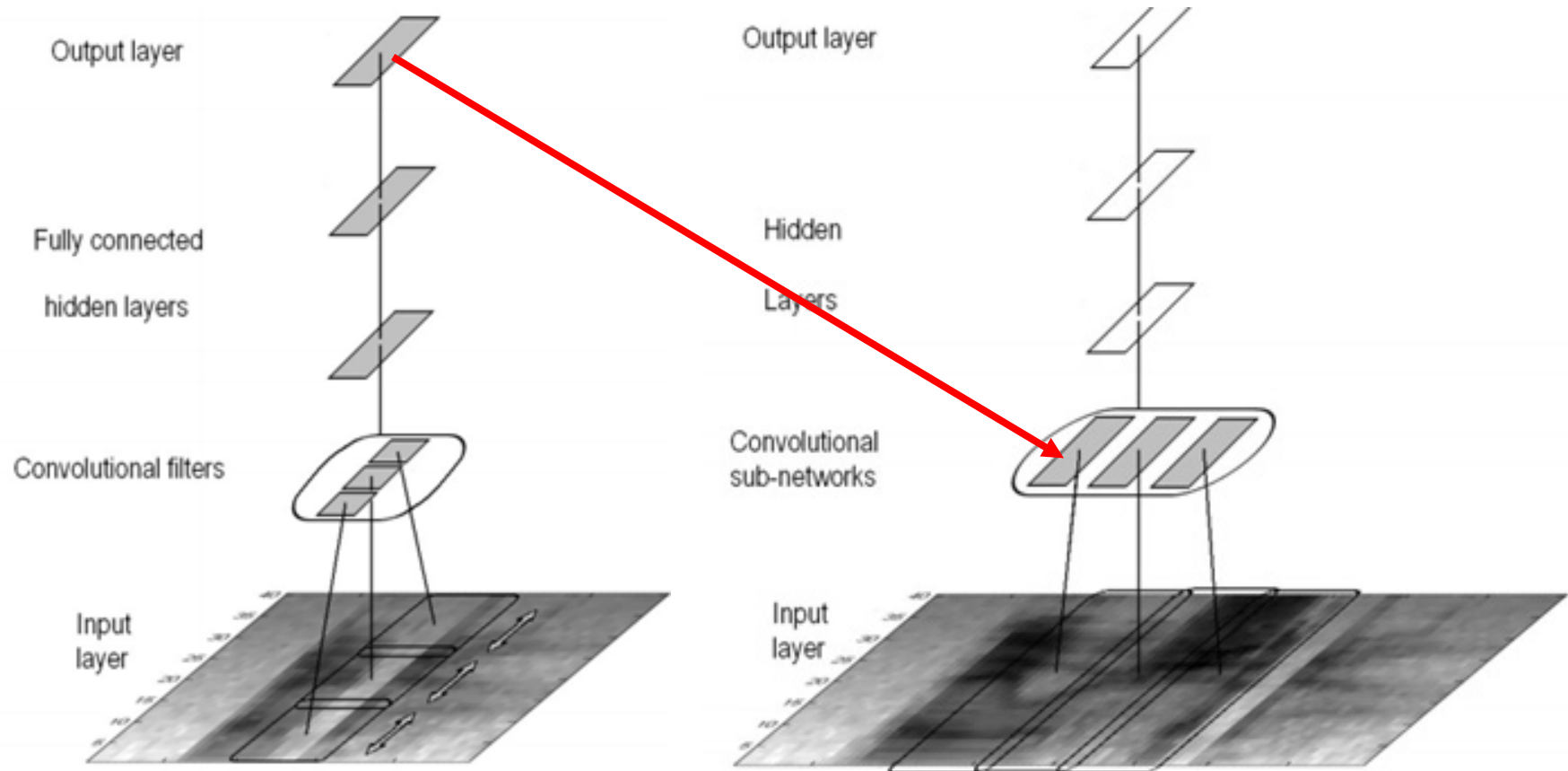
Probabilities of states



Replace DNN by CNN

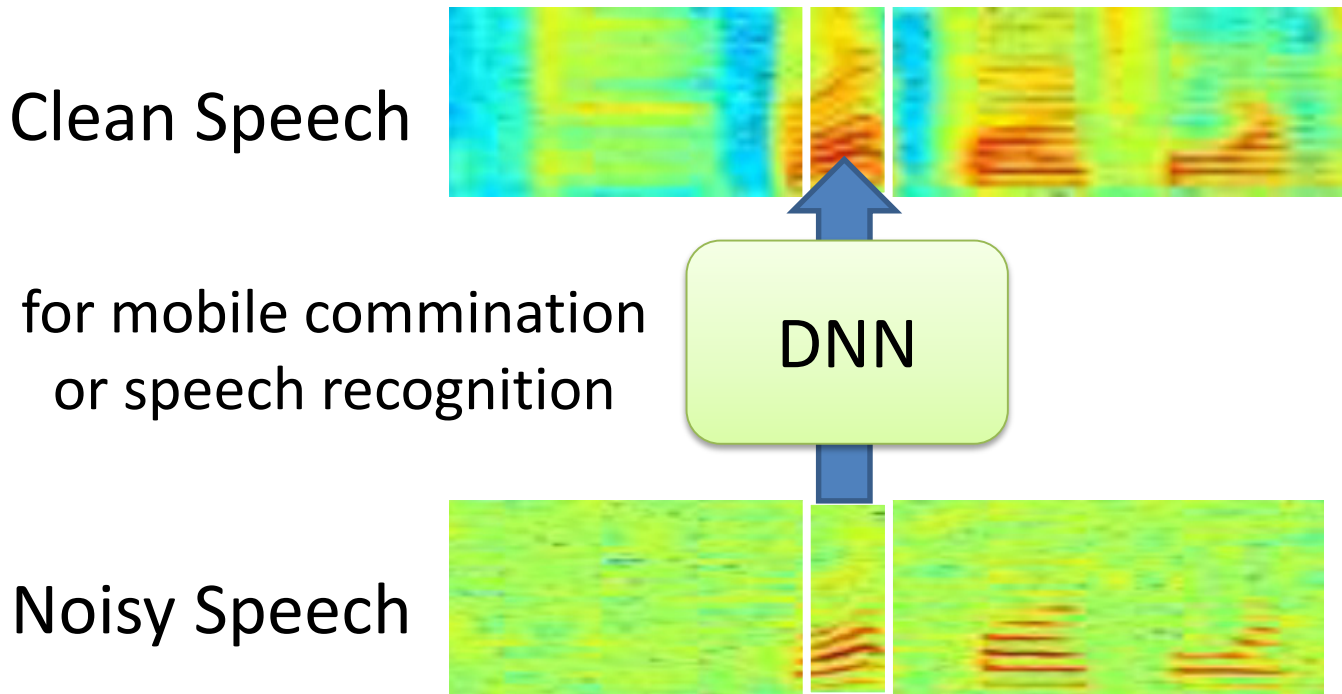






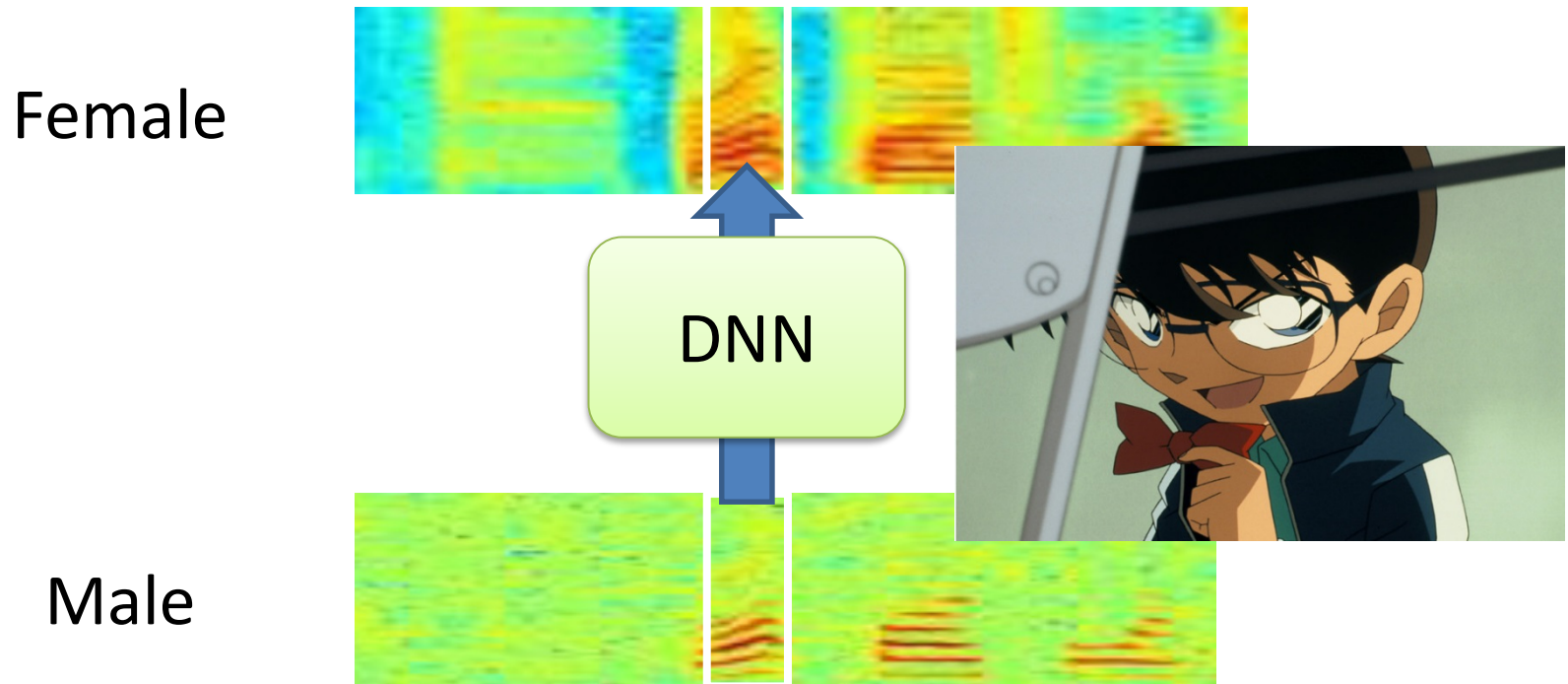
Applications in Acoustic Signal Processing





- Demo for speech enhancement:

http://home.ustc.edu.cn/~xuyong62/demo/SE_DNN.html



- Demo for Voice Conversion:
<https://candyvoice.com/demos/voice-conversion?lang=en>

- ◆ google IO 2018 assistant call DEMO1 (彩色=AI)



- ◆ google IO 2018 assistant call DEMO2 (彩色=AI)

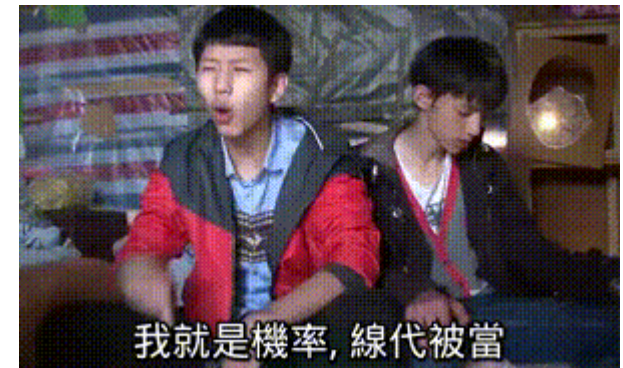


Concluding Remarks



It's an interesting time!

- ◆ 大學的基礎科目很重要！
 - 機率統計, 線性代數, Machine Learning, ... etc
- ◆ Deep learning integrated into standard speech toolkits
 - Kaldi, HTK, ... etc
- ◆ Rich variety of models and topologies supported by:
 - large quantities of training data
 - GPU-based training (and parallel implementations)
 - array of **ML tools**: TensorFlow, CNTK... etc



Questions

Thank you

To learn more about Delta, please visit www.deltaww.com.



- [1] L. Baum and J. Eagon, “An Inequality with Applications to Statistical Estimation for Probabilistic Functions of Markov Processes and to a Model for Ecology,” *Bull Amer Math Soc*, vol. 73, pp. 360–363, 1967.
- [2] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer Verlag, 2006.
- [3] H. Bourlard and N. Morgan, “Connectionist speech recognition: A hybrid approach,” 1994.
- [4] W. Chan, N. Jaitly, Q. V. Le, and O. Vinyals, “Listen, attend and spell,” *CoRR*, vol. abs/1508.01211, 2015.
[Online]. Available: <http://arxiv.org/abs/1508.01211>
- [5] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, “Attention-based models for speech recognition,” *CoRR*, vol. abs/1506.07503, 2015.
[Online]. Available: <http://arxiv.org/abs/1506.07503>
- [6] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, “Empirical evaluation of gated recurrent neural networks on sequence modeling,” *arXiv preprint arXiv:1412.3555*, 2014.
- [7] J. Chung, K. Kastner, L. Dinh, K. Goel, A. C. Courville, and Y. Bengio, “A recurrent latent variable model for sequential data,” *CoRR*, vol. abs/1506.02216, 2015.
[Online]. Available: <http://arxiv.org/abs/1506.02216>
- [8] M. Gales and S. Young, “The application of hidden Markov models in speech recognition,” *Foundations and Trends in Signal Processing*, vol. 1, no. 3, 2007.
- [9] A. Graves, “Sequence transduction with recurrent neural networks,” *CoRR*, vol. abs/1211.3711, 2012.
[Online]. Available: <http://arxiv.org/abs/1211.3711>
- [10] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd international conference on Machine learning*. ACM, 2006, pp. 369–376.
- [11] A. Graves, A.-R. Mohamed, and G. Hinton, “Speech recognition with deep recurrent neural networks,” in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013, pp. 6645–6649.
- [12] G. E. Hinton, “Products of experts,” in *Proceedings of the Ninth International Conference on Artificial Neural Networks (ICANN 99)*, 1999, pp. 1–6. 57/57

- [13] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-R. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, et al., “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [14] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [15] F. Jelinek, *Statistical methods for speech recognition*, ser. Language, speech, and communication. Cambridge (Mass.), London: MIT Press, 1997.
- [16] H.-K. Kuo and Y. Gao, “Maximum entropy direct models for speech recognition,” *IEEE Transactions Audio Speech and Language Processing*, 2006.
- [17] L. Lu, X. Zhang, K. Cho, and S. Renals, “A study of the recurrent neural network encoder-decoder for large vocabulary speech recognition,” in *Proc. INTERSPEECH*, 2015.
- [18] T. Robinson and F. Fallside, “A recurrent error propagation network speech recognition system,” *Computer Speech & Language*, vol. 5, no. 3, pp. 259–274, 1991.
- [19] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1,” D. E. Rumelhart, J. L. McClelland, and C. PDP Research Group, Eds. Cambridge, MA, USA: MIT Press, 1986, ch. Learning Internal Representations by Error Propagation, pp. 318–362.
- [20] T. N. Sainath, O. Vinyals, A. Senior, and H. Sak, “Convolutional, long short-term memory, fully connected deep neural networks,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 4580–4584.
- [21] M. Schuster and K. K. Paliwal, “Bidirectional recurrent neural networks,” *IEEE Transactions on Signal Processing*, vol. 45, no. 11, pp. 2673–2681, 1997.
- [22] R. K. Srivastava, K. Greff, and J. Schmidhuber, “Highway networks,” *CoRR*, vol. abs/1505.00387, 2015.
[Online]. Available: <http://arxiv.org/abs/1505.00387>
- [23] A. Viterbi, “Error bounds for convolutional codes and an asymptotically optimum decoding algorithm,” vol. 13, no. 2, pp. 260–269, 1967.
- [24] Research Talk by Mark Gales, in *LXMLS*, 2017. <http://lxmls.it.pt/2017/talk.pdf>
- [25] NTU Course by Hung-Yee Lee ([ppt](#))
- [26] Stanford Course: Speech Processing <https://web.stanford.edu/class/cs224s/lectures/>

- ◆ **Speech-to-text (Automatic Speech Recognition)**
- ◆ **Text -to-speech (Speech Synthesis)**
- ◆ **Speaker Recognition & Diarization**
- ◆ **Voice Conversion**
- ◆ **Speech Separation**
- ◆ **Speech Enhancement**
- ◆ **Speech Emotion Recognition**
- ◆ **Music Information Retrieval**
- ◆ **Spoken Language Understanding**
- ◆ ...