

Spoken Content Retrieval

Beyond Cascading Speech Recognition and Text Retrieval

Lin-shan Lee and Hung-yi Lee
National Taiwan University

Focus of this Tutorial



- New frontiers and directions towards the future of speech technologies
- Not skills and experiences in optimizing performance in evaluation programs

Text Content Retrieval



Obama



Google 搜尋

好手氣

Voice Search

約有 140,000,000 項結果 (搜尋時間：0.25 秒)

[Barack Obama](#)

[www.barackobama.com/](#) ▾ 翻譯這個網頁

Official re-election campaign website of President Barack Obama provides the latest updates, election news, videos, local events and ways to volunteer and ...

[Barack Obama - Wikipedia, the free encyclopedia](#)

[en.wikipedia.org/wiki/Barack_Obama](#) ▾

In 2004, Obama received national attention during his campaign to represent Illinois in the United States Senate with his victory in the March Democratic Party ...

Family of Barack Obama - Early life and career of Barack - Michelle Obama - Luo

[President Barack Obama | The White House](#)

[www.whitehouse.gov](#) ▸ The Administration ▾ 翻譯這個網頁

Spoken Content Retrieval



Obama



Google 搜尋

好手氣

Spoken Content



Lectures



Broadcast Program

Multimedia Content

約有 140,000,000 項結果 (搜尋時間: 0.25 秒)

[Barack Obama](#)

www.barackobama.com/ ▾ 翻譯這個網頁

Official re-election campaign website of President Barack Obama. Updates, election news, videos, local events and ways to get involved.

[Barack Obama - Wikipedia, the free encyclopedia](#)
en.wikipedia.org/wiki/Barack_Obama ▾

In 2004, Obama received national attention during his victory in the United States Senate with his victory in the March 2004 election. Family of Barack Obama - Early life and career of Barack Obama.

[President Barack Obama | The White House](#)

www.whitehouse.gov ▾ The Administration ▾ 翻譯這個網頁

Spoken Content Retrieval



300 hrs multimedia is
uploaded per minute.
(2015.01)

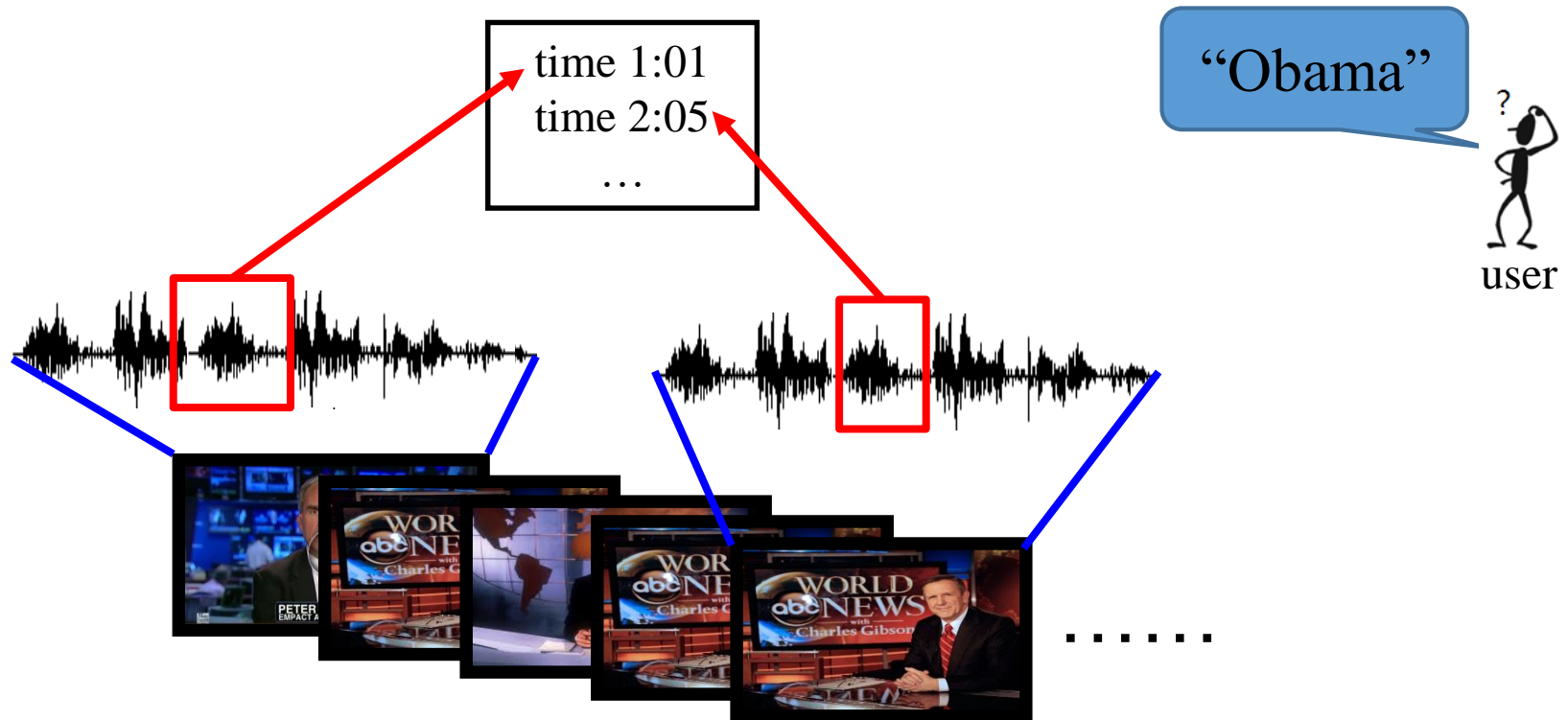


1874 courses on coursera
(2016.04)

- Nobody is able to go through the data.
- In these multimedia, the spoken part carries very important information about the content
- Spoken content retrieval: Machine listens to the data, and extract the desired information for each individual user.
 - Just as Google does on text data

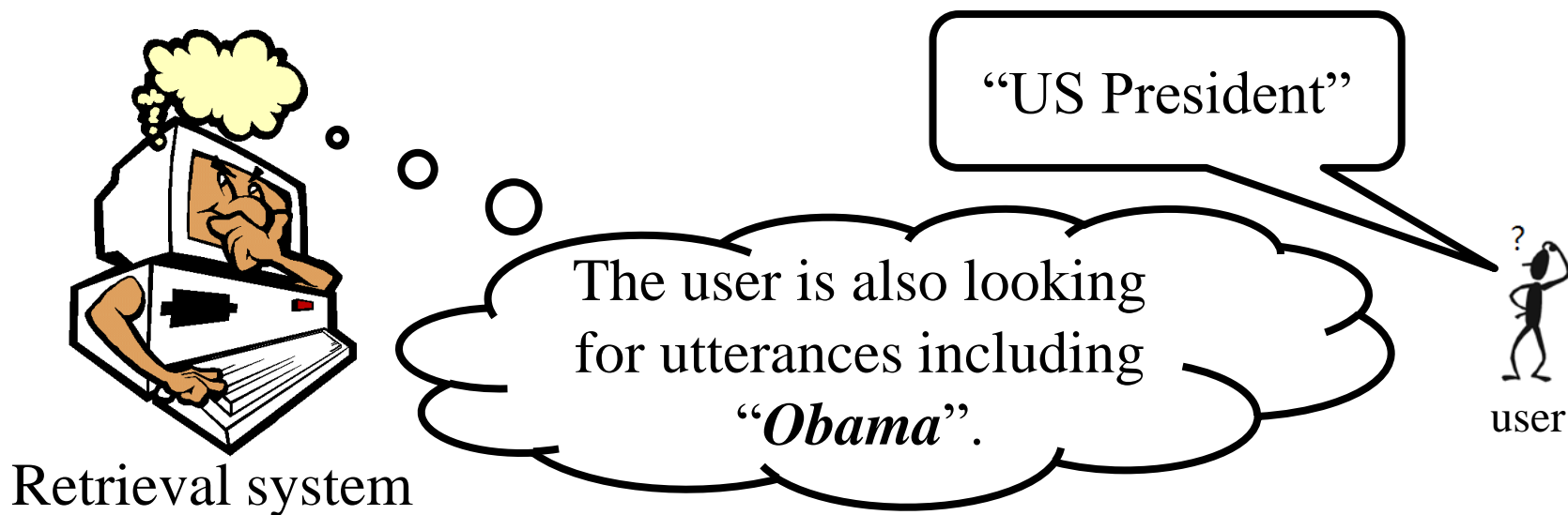
Spoken Content Retrieval – Goal

- Basic goal: Identify the time spans that the query occurs in an audio database
 - This is called “*Spoken Term Detection*”



Spoken Content Retrieval – Goal

- Basic goal: Identify the time spans that the query occurs in an audio database
 - This is called “*Spoken Term Detection*”
- Advanced goal: Semantic retrieval of spoken content

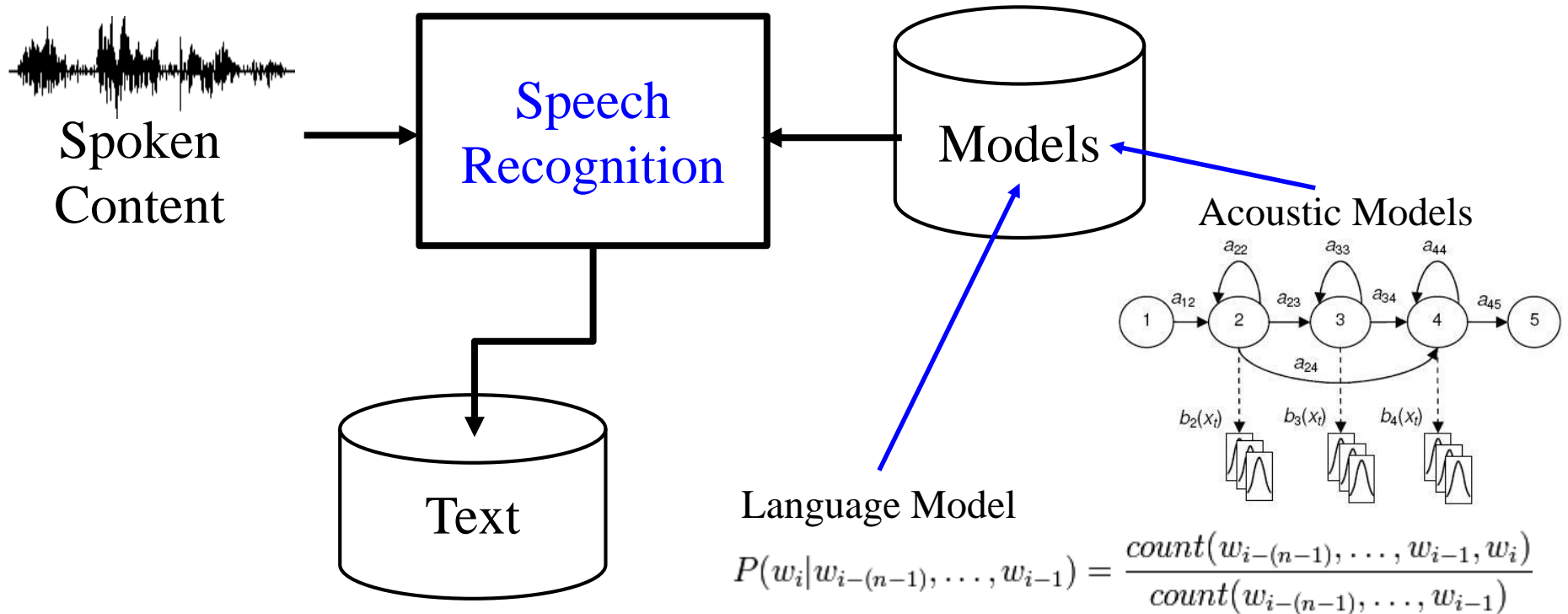


It is natural to think



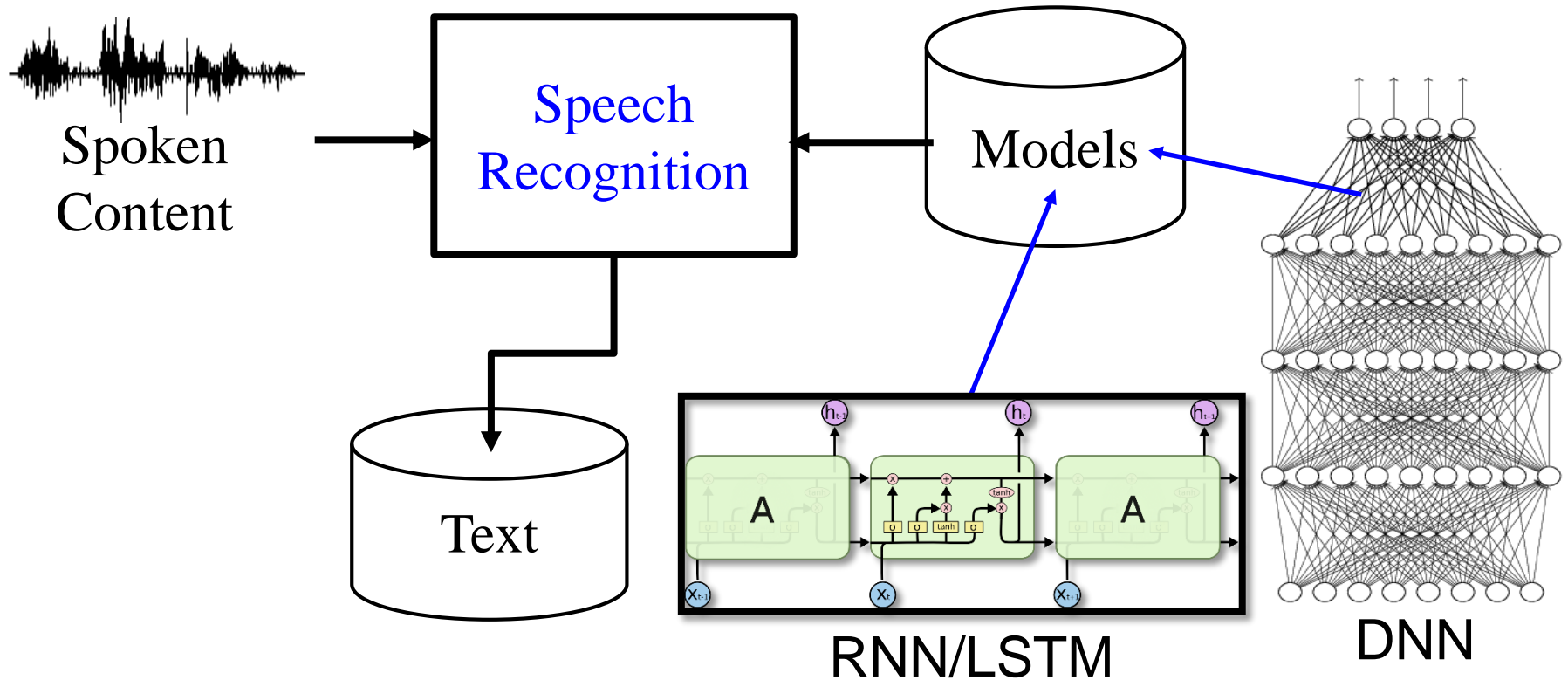
Spoken Content Retrieval
||
Speech Recognition
+
Text Retrieval

It is natural to think



- Transcribe spoken content into text by speech recognition

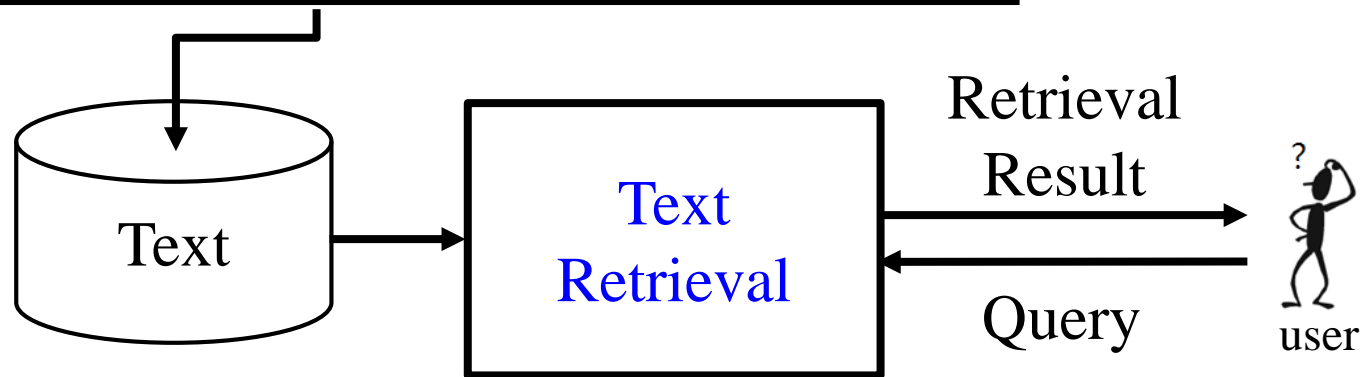
It is natural to think



- Transcribe spoken content into text by speech recognition

It is natural to think

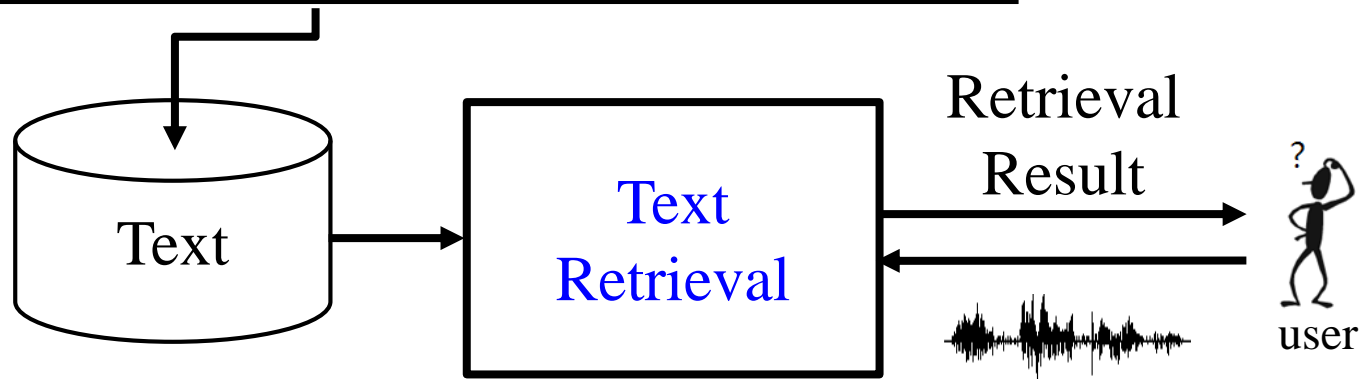
Black Box



- Transcribe spoken content into text by speech recognition
- Use text retrieval approaches to search over the transcriptions

It is natural to think

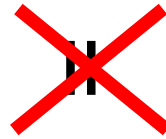
Black Box



- For spoken queries, transcribe them into text by speech recognition.

Our point in this tutorial

Spoken Content Retrieval



Speech Recognition

+

Text Retrieval

Outline

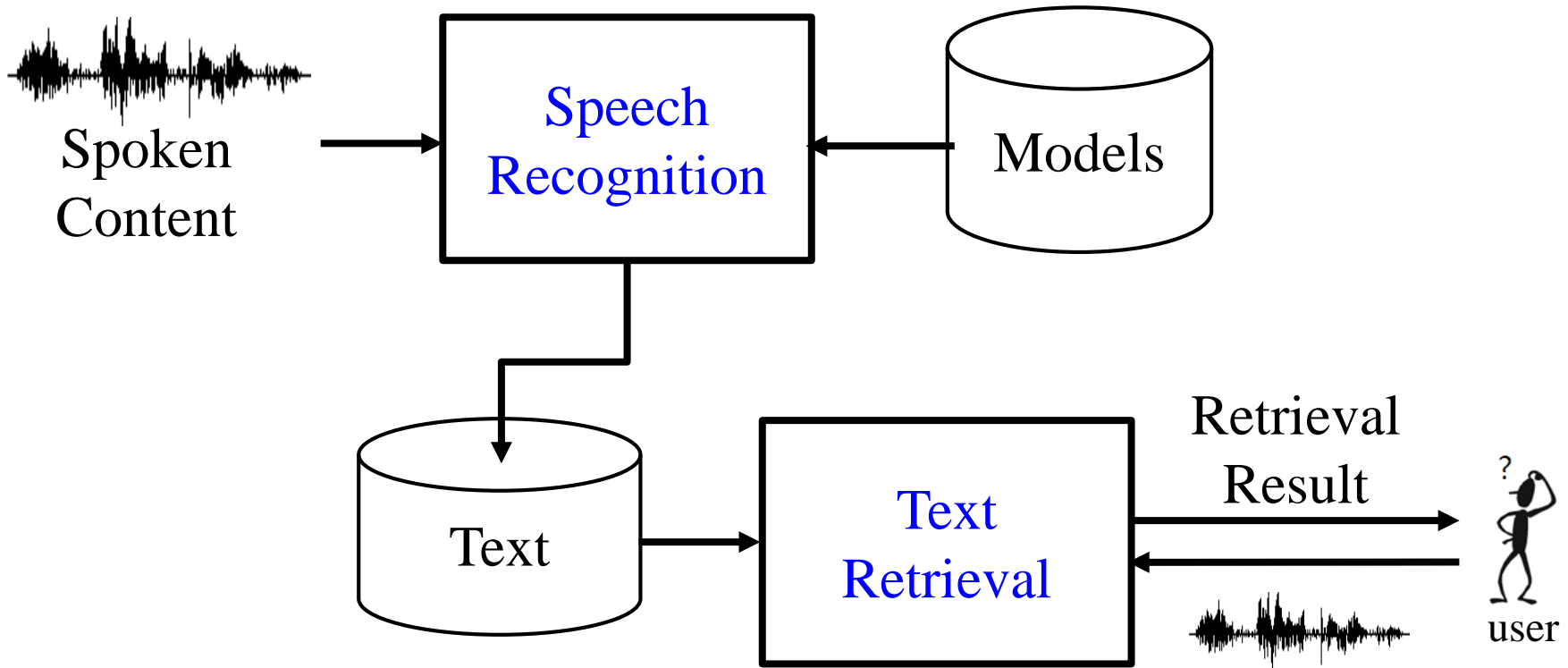
- Introduction: Conventional Approach:
*Spoken Content Retrieval =
Speech Recognition + Text Retrieval*
- Core: Beyond Cascading Speech
Recognition and Text Retrieval
 - ▣ Five new directions

Introduction:

Spoken Content Retrieval =

Speech Recognition + Text Retrieval

It is natural to think

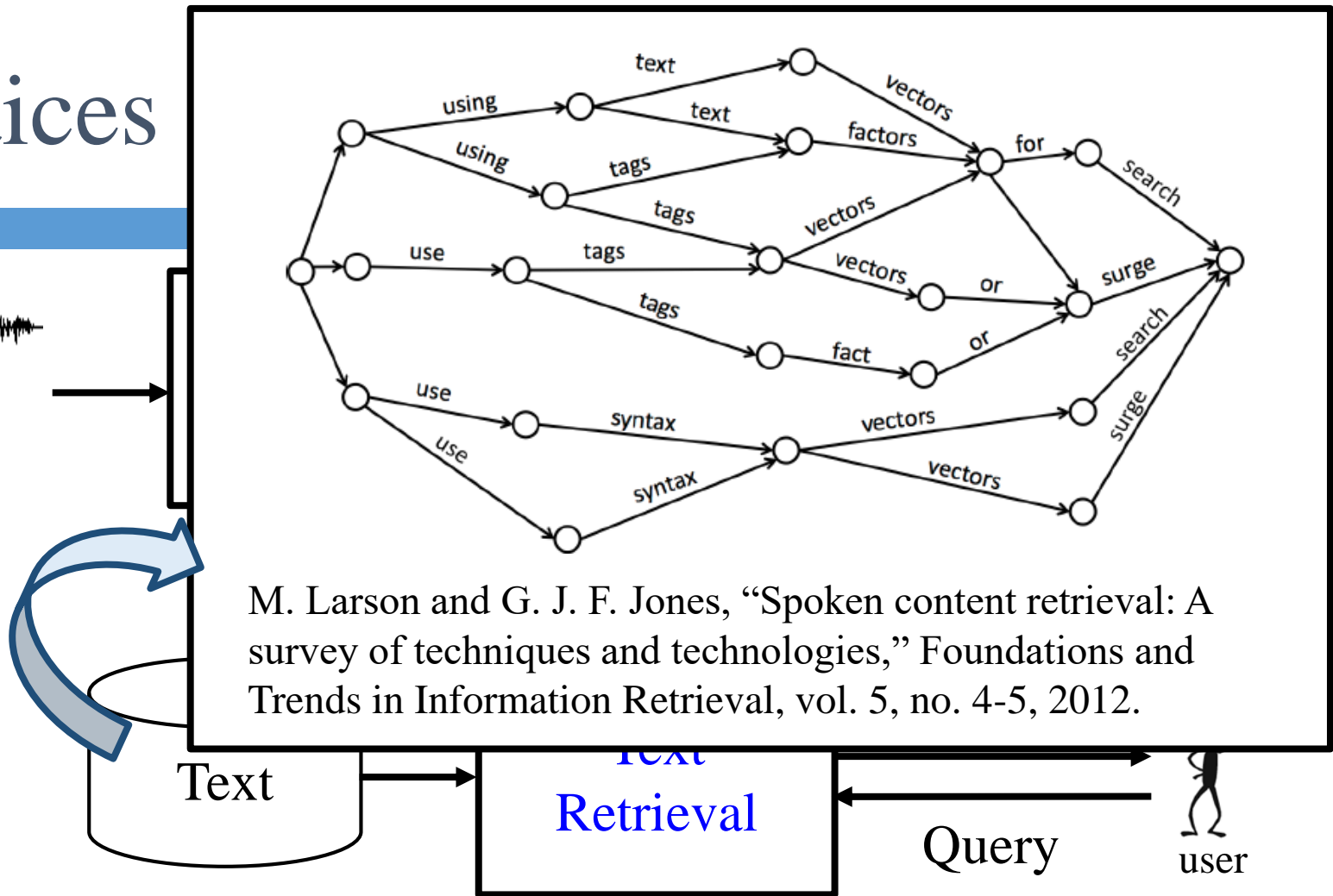


Speech Recognition always produces errors.

Lattices



Spoken
Content

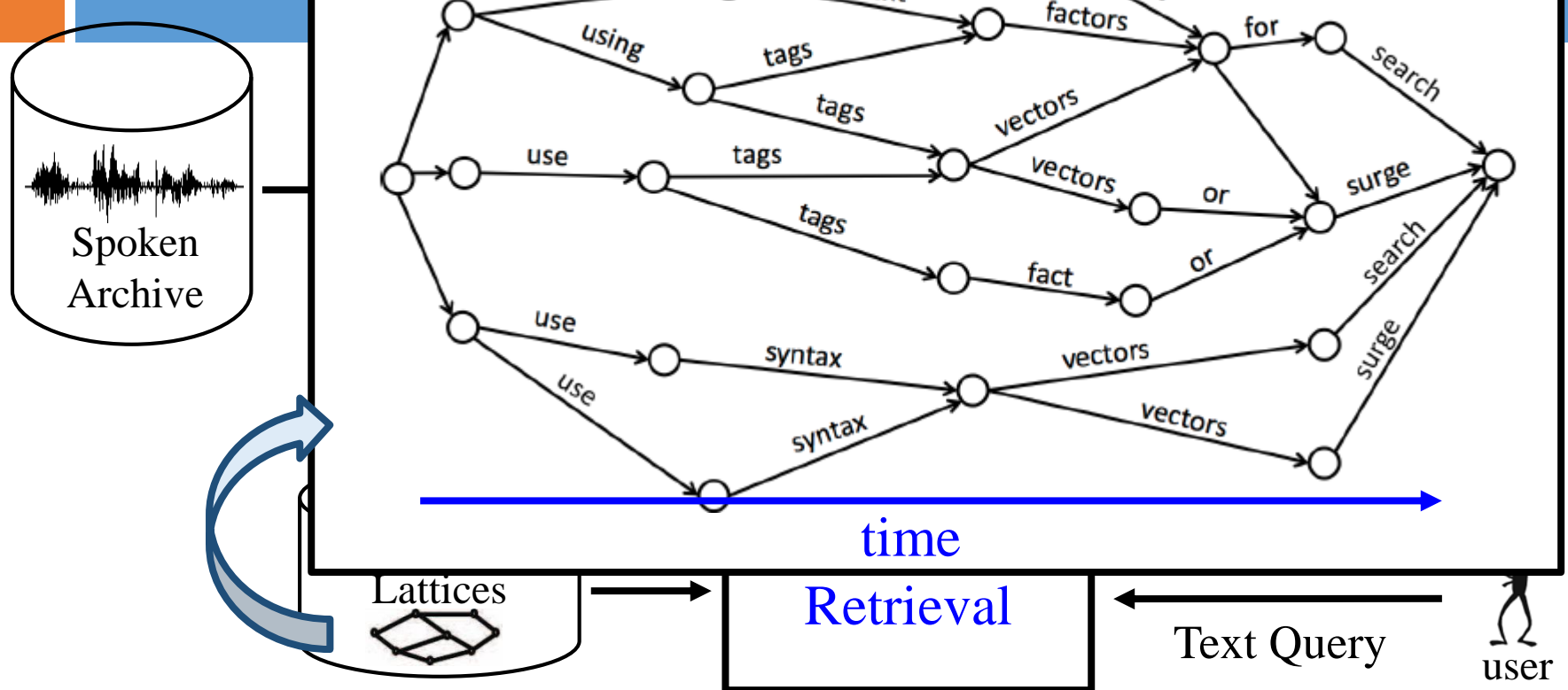


Lattices



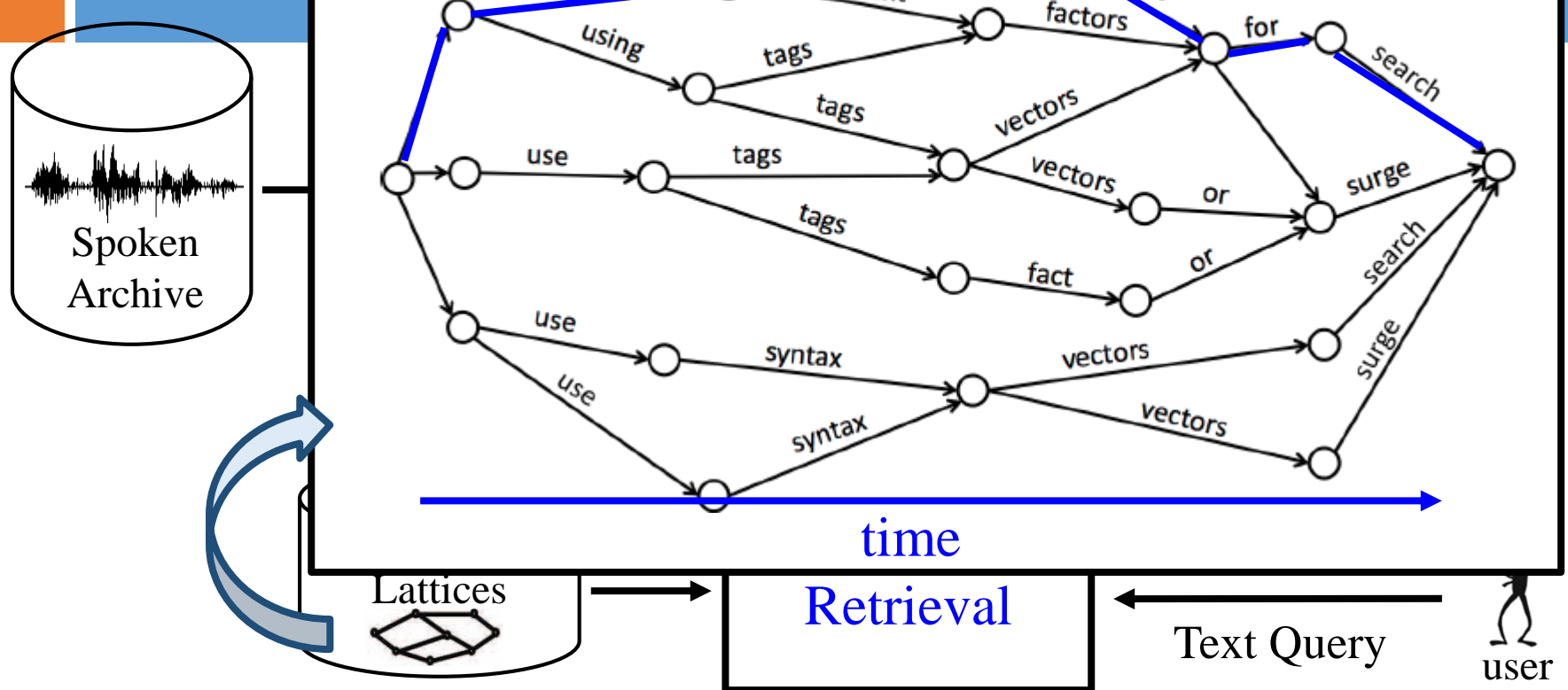
- Keep most possible recognition output
- Each path has a weight (confidence to be correct)

Lattice



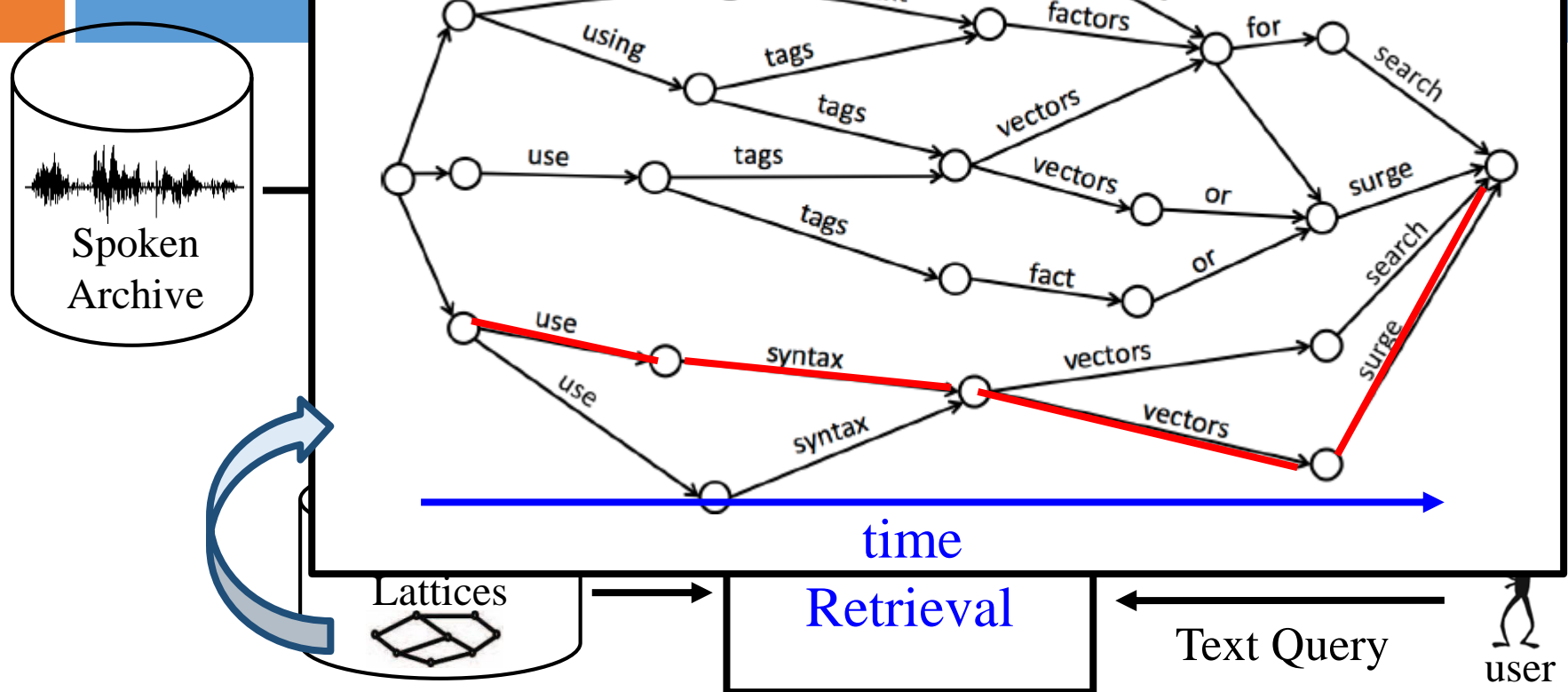
- Horizontal scale is the time
- Each path is a possible recognition result

Lattice



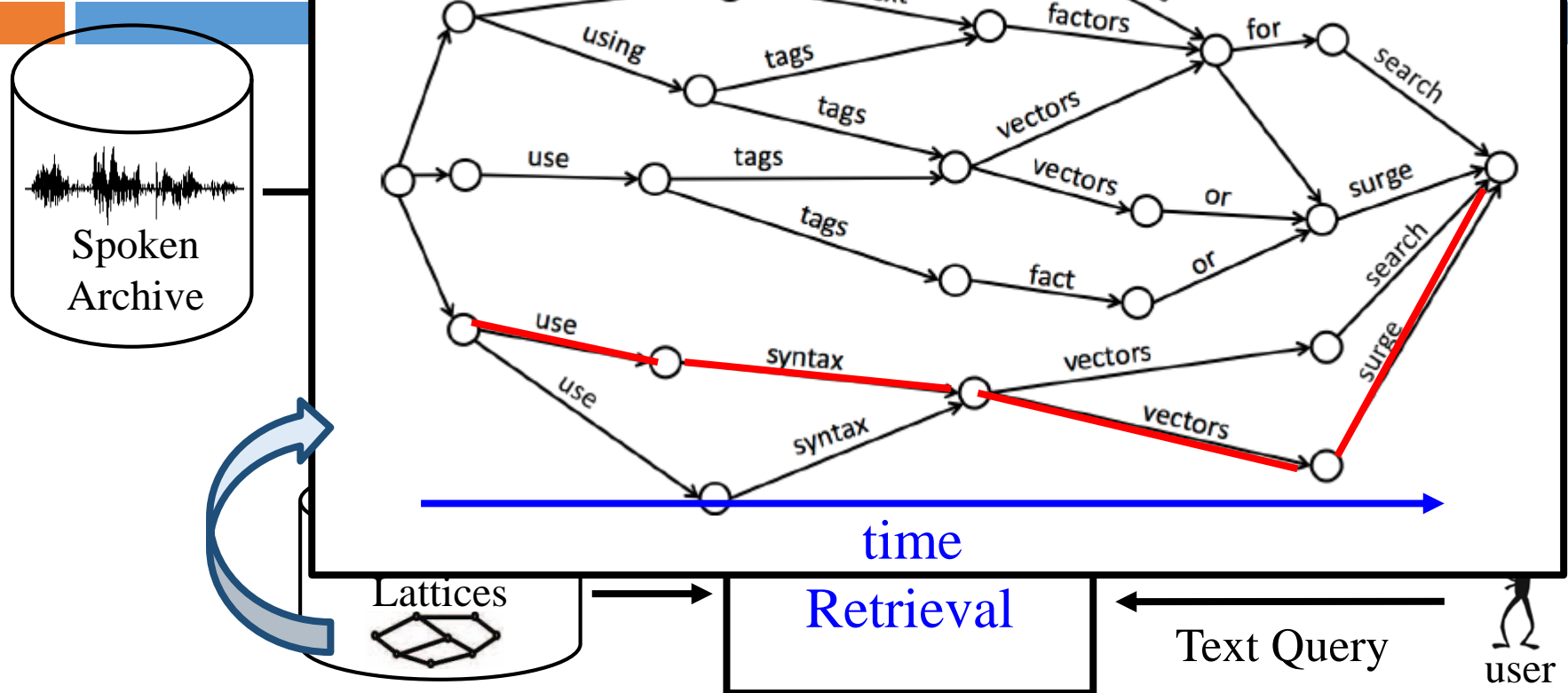
- Horizontal scale is the time
- Each path is a possible recognition result

Lattice



- Horizontal scale is the time
- Each path is a possible recognition result

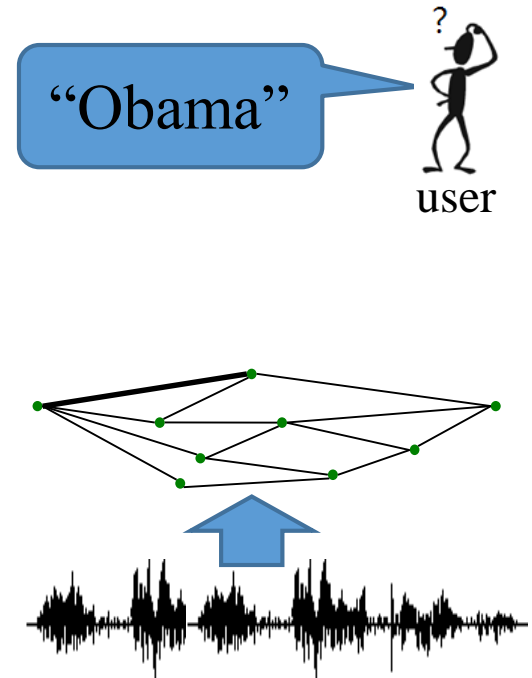
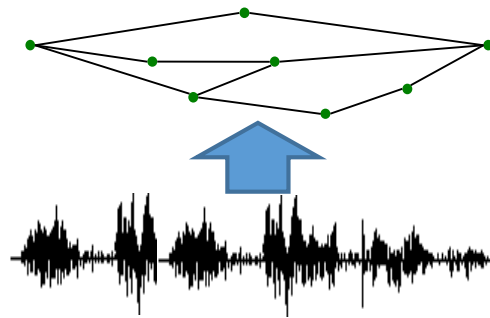
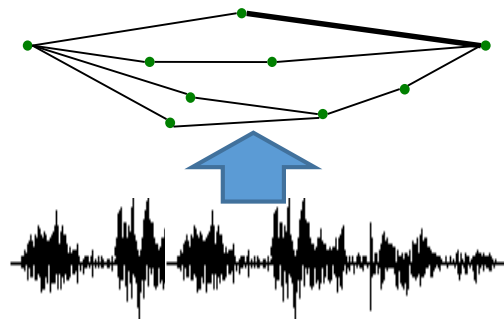
Lattice



- Higher probability to include the correct words
- More noisy words included inevitably
- Higher memory/computation requirements

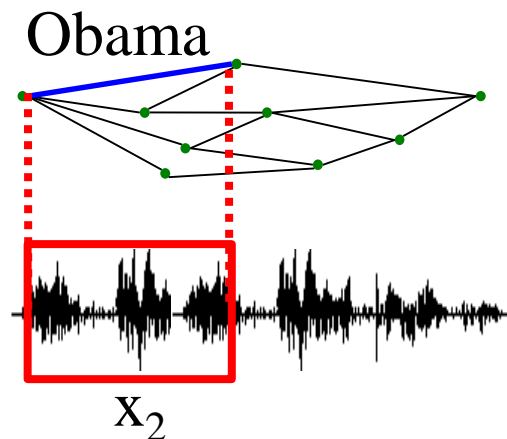
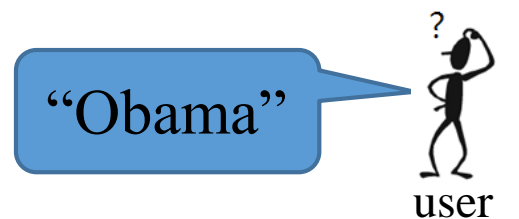
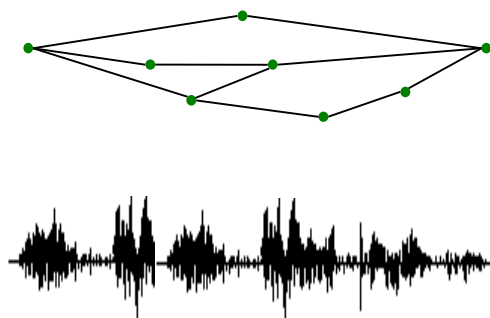
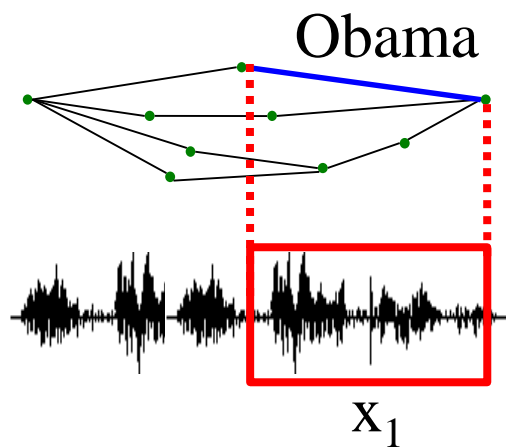
Searching over Lattices

- Consider the basic goal: Spoken Term Detection



Searching over Lattices

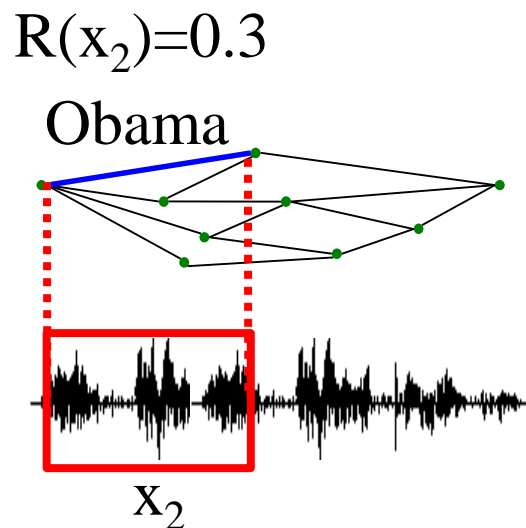
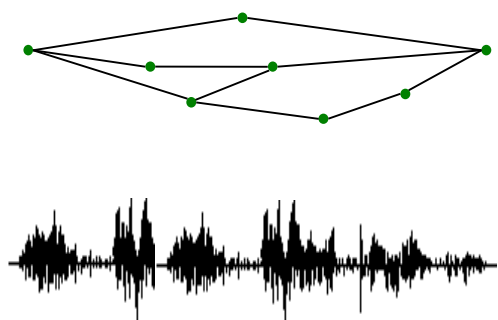
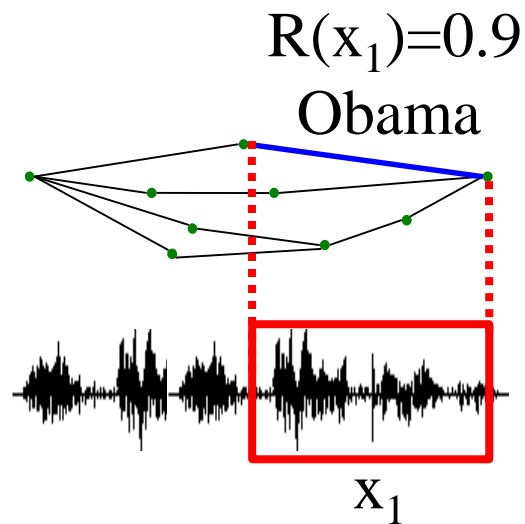
- Consider the basic goal: Spoken Term Detection
 - ▣ Find the arcs hypothesized to be the query term



Searching over Lattices

- Consider the basic goal: Spoken Term Detection
 - ▣ *Posterior probabilities* computed from lattices used as confidence scores

Two ways to display the results:
unranked and **ranked**.

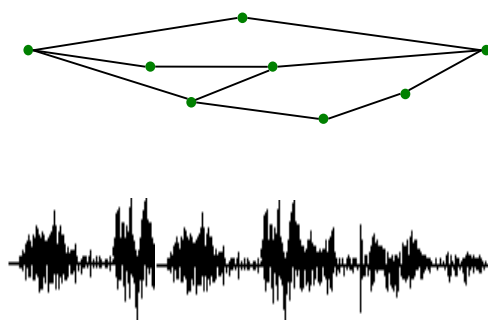
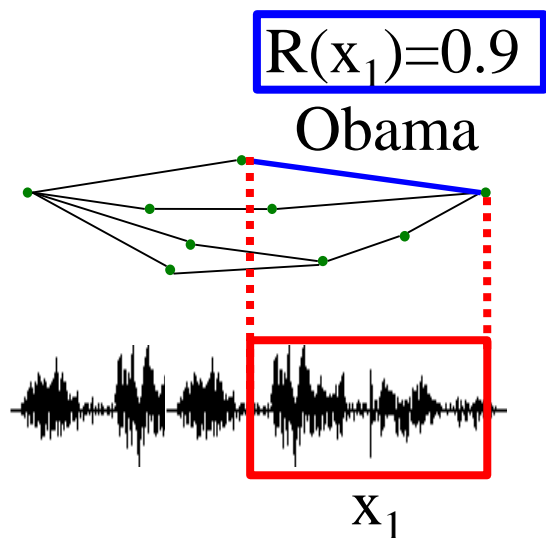


Searching over Lattices

- Consider the basic goal: Spoken Term Detection
 - ▣ **Unranked**: Return the results with the scores higher than a threshold

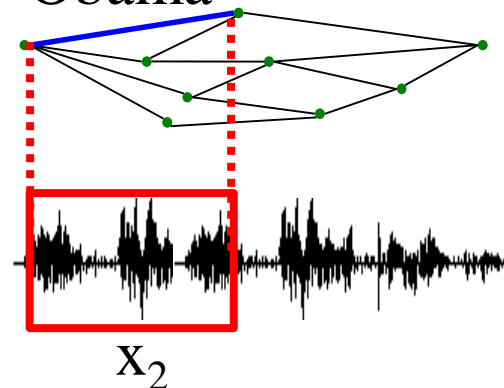
Set the threshold as 0.6

Return x_1



$R(x_2)=0.3$

Obama

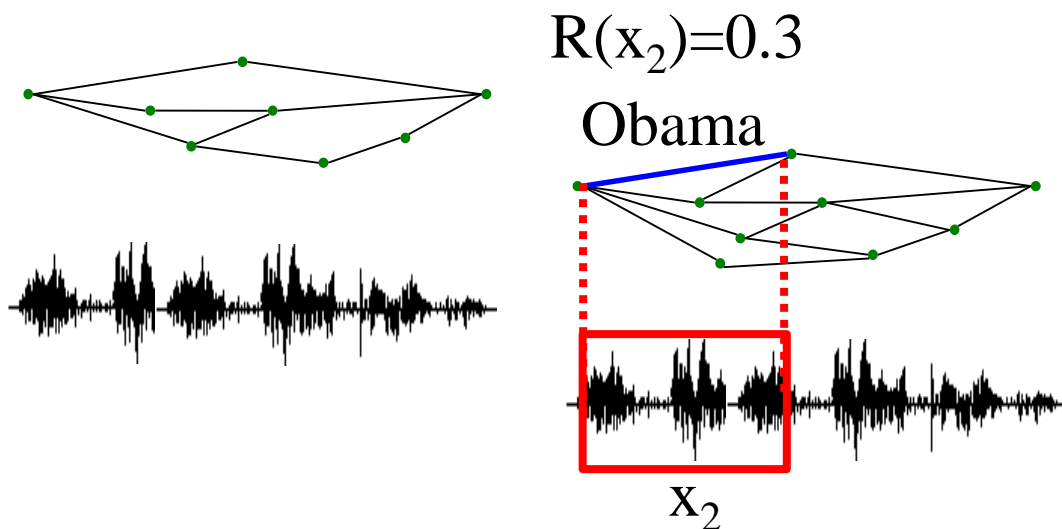
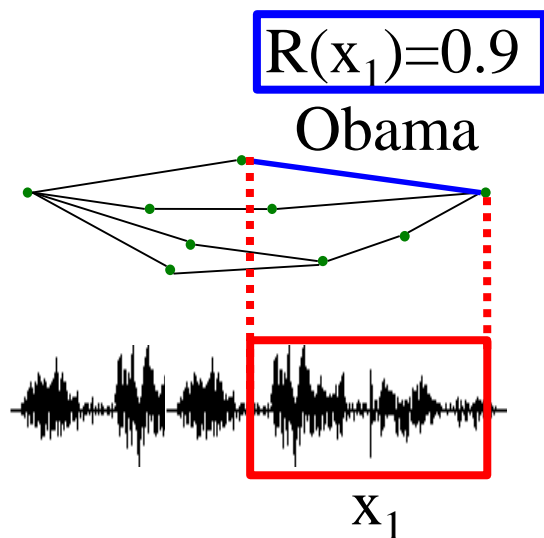


Searching over Lattices

- Consider the basic goal: Spoken Term Detection
 - ▣ **Unranked**: Return the results with the scores higher than a threshold

The threshold can be determined automatically and query specific.


[Miller, Interspeech 07][Can, HLT 09][Mamou, ICASSP 13][Karakos, ASRU 13][Zhang, Interspeech 12][Pham, ICASSP 14]



Actual Term Weighted Value (ATWV)

□ Evaluating unranked result

retrieved



threshold

time 1:01	1.0
time 2:05	0.9
time 1:31	0.7
.....	

The table is enclosed in a blue rounded rectangle. A red horizontal line is drawn across the table, separating the top two rows from the bottom two rows. The word 'retrieved' is positioned above the top two rows, and 'threshold' is positioned below the red line.

$$ATWV = 1 - P_{miss} - \beta P_{FA}$$

$$P_{miss} = 1 - N_{correct} / N_{ref}$$

$$P_{FA} = N_{spurious} / N_{NT}$$

N_{ref} : number of times the query term appears in audio database

$N_{correct}$: the number of retrieved objects that are actually **correct**

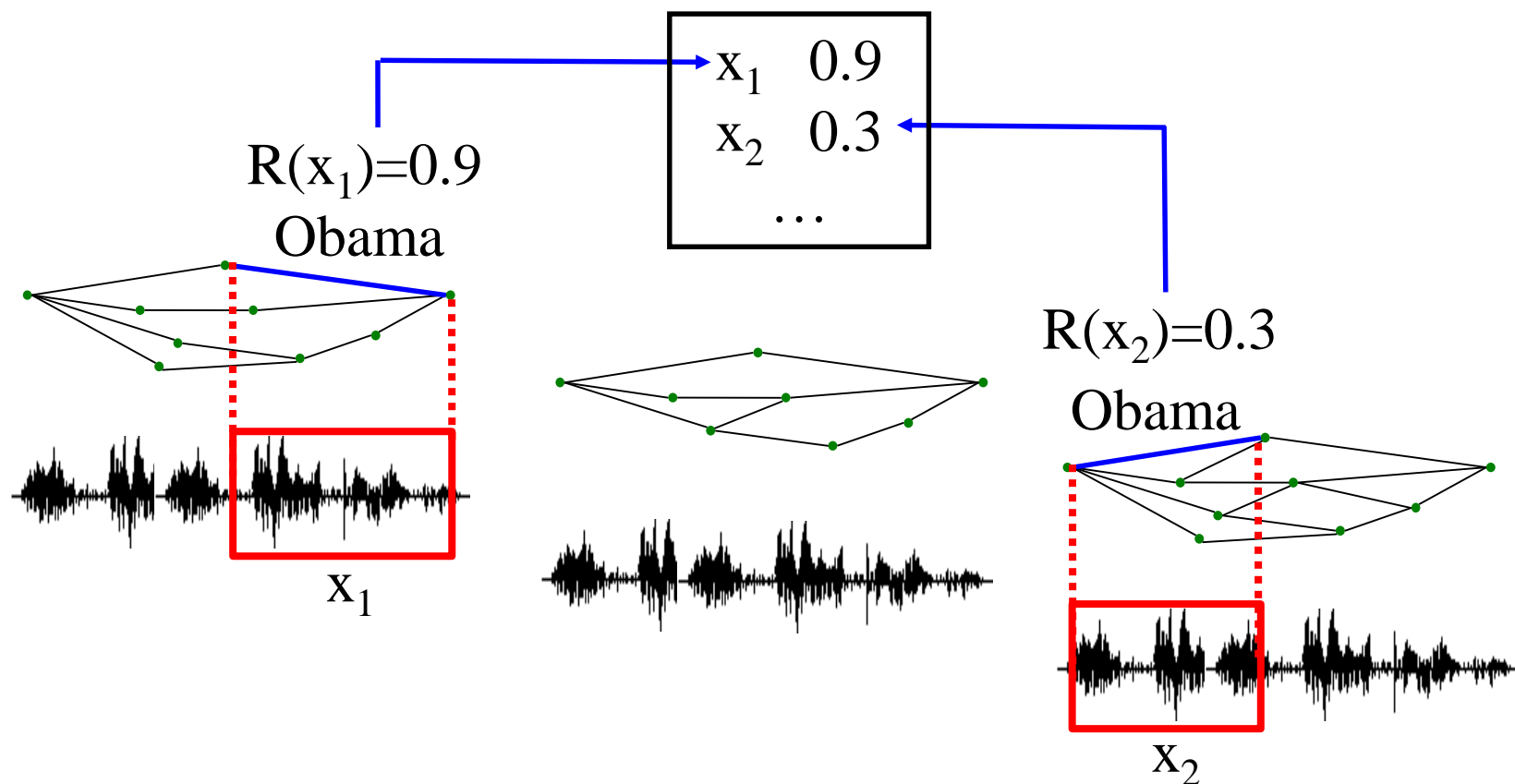
$N_{spurious}$: the number of retrieved objects that are **incorrect**

N_{NT} : audio duration (in seconds) – N_{ref}

Maximum Term Weighted Value (MTWV): tune the threshold to obtain the best ATWV

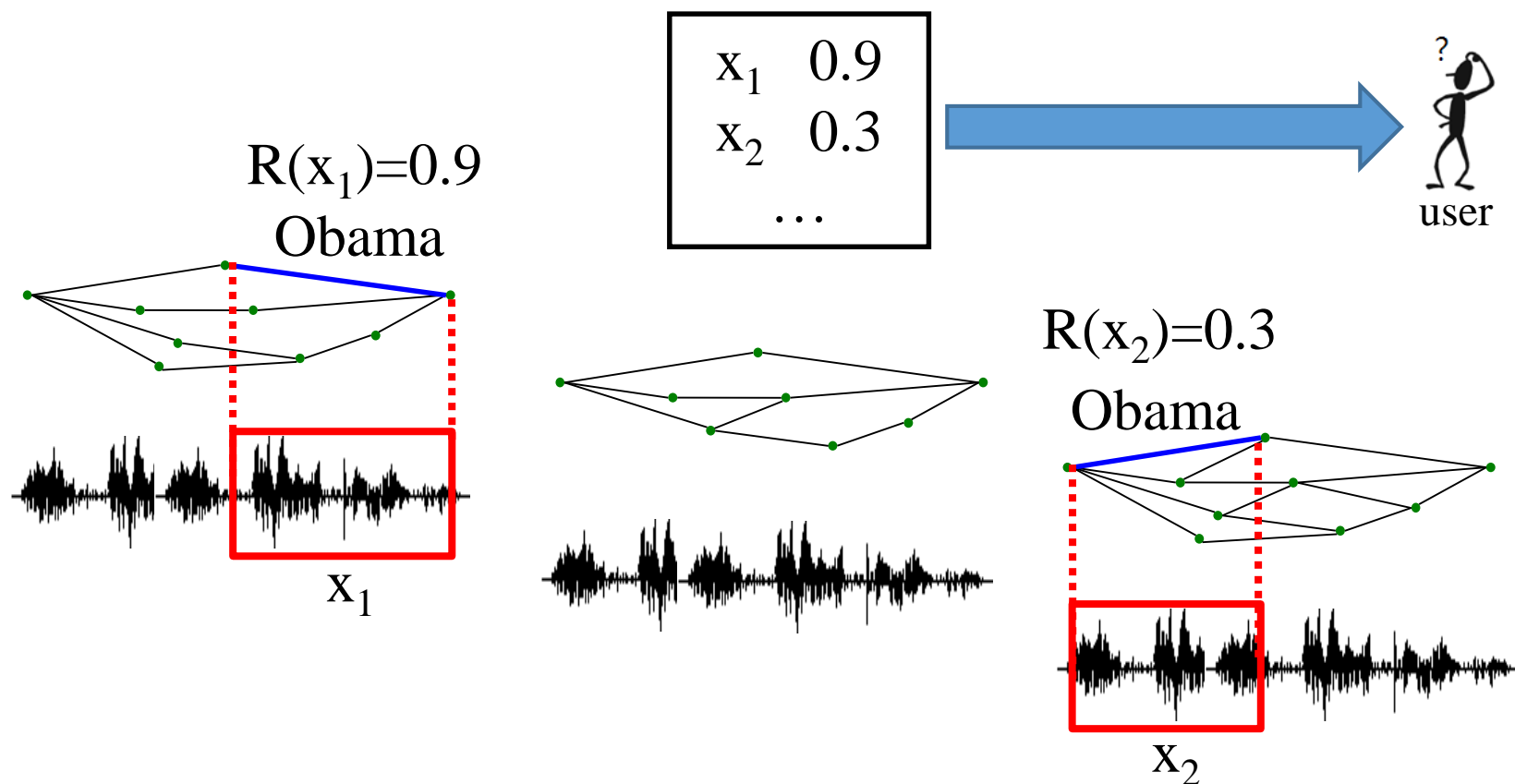
Searching over Lattices

- Consider the basic goal: Spoken Term Detection
 - ▣ **Ranked:** results ranked according to the scores



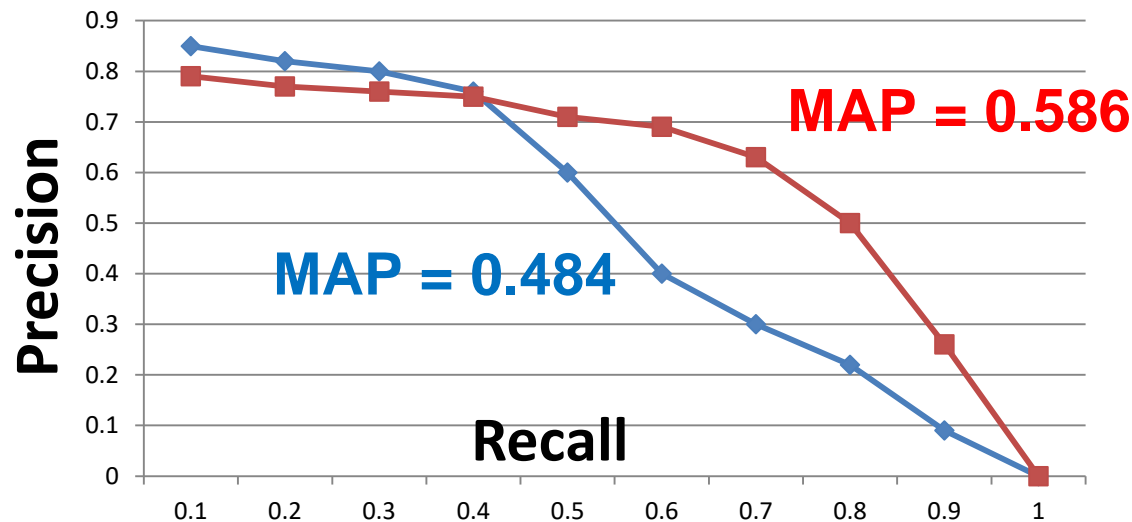
Searching on Lattices

- Consider the basic goal: Spoken Term Detection
 - ▣ **Ranked:** The results are ranked according to the scores



Mean Average Precision (MAP)

- Evaluating ranked list
- area under recall-precision curve
 - Recall: percentage of ground truth results retrieved
 - Precision: percentage of retrieved results being correct
 - Higher threshold gives higher precision but lower recall, etc.



Examples of Lattice Indexing

Approaches

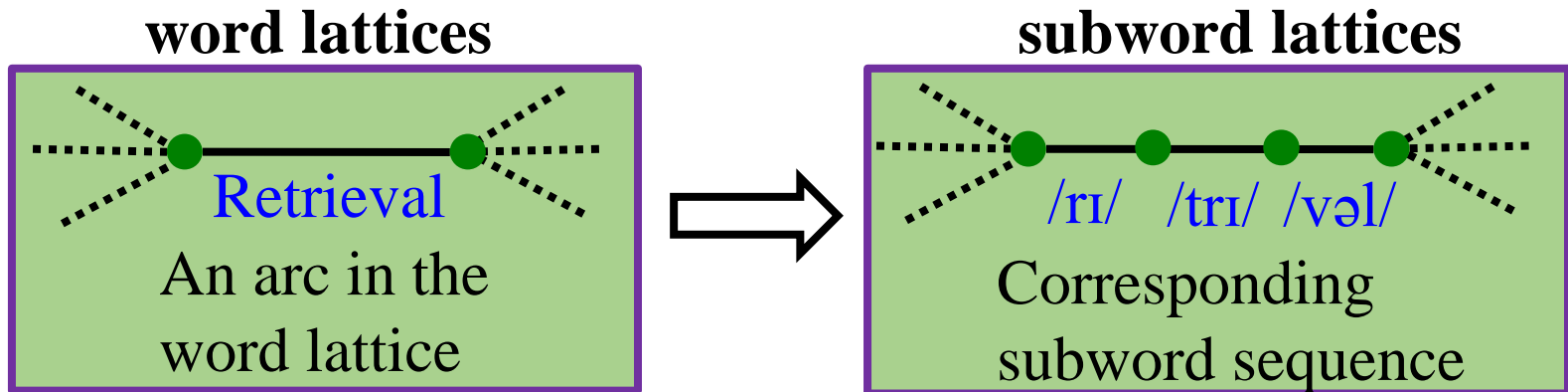
- Position Specific Posterior Lattices (PSPL)[Chelba, ACL 05][Chelba, Computer Speech and Language 07]
- Confusion Networks (CN)[Mamou, SIGIR 06][Hori, ICASSP 07][Mamou, SIGIR 07]
- Time-based Merging for Indexing (TMI)[Zhou, HLT 06][Seide, ASRU 07]
- Time-anchored Lattice Expansion (TALE)[Seide, ASRU 07][Seide, ICASSP 08]
- WFST: directly compile the lattice into a weighted finite state transducer [Allauzen, HLT 04][Parlak, ICASSP 08][Can, ICASSP 09][Parada, ASRU 09]

Out-of-Vocabulary (OOV) Problem

- Speech recognition is based on a lexicon
- Words not in the lexicon can never be transcribed
- Many informative words are out-of-vocabulary (OOV)
 - Many query terms are new or special words or named entities

Subword-based Retrieval

- All OOV words composed of subword units
 - ▣ Generate subword lattices
 - Transform word lattices into subword lattices



- Can also be directly generated by speech recognition using subword-based lexicon and language model

Subword-based Retrieval

- Subword-based retrieval
 - ▣ Generate subword lattices
 - ▣ Transform user query into subword sequence
 - Obama → /au/ /ba/ /mə/
 - ▣ Text retrieval techniques equally useful except based on subword lattices and subword query
 - Replace words by subword units
 - ▣ OOV words can be retrieved by matching over the subword units without being recognized

Subword-based Retrieval

- Frequently Used Subword Units

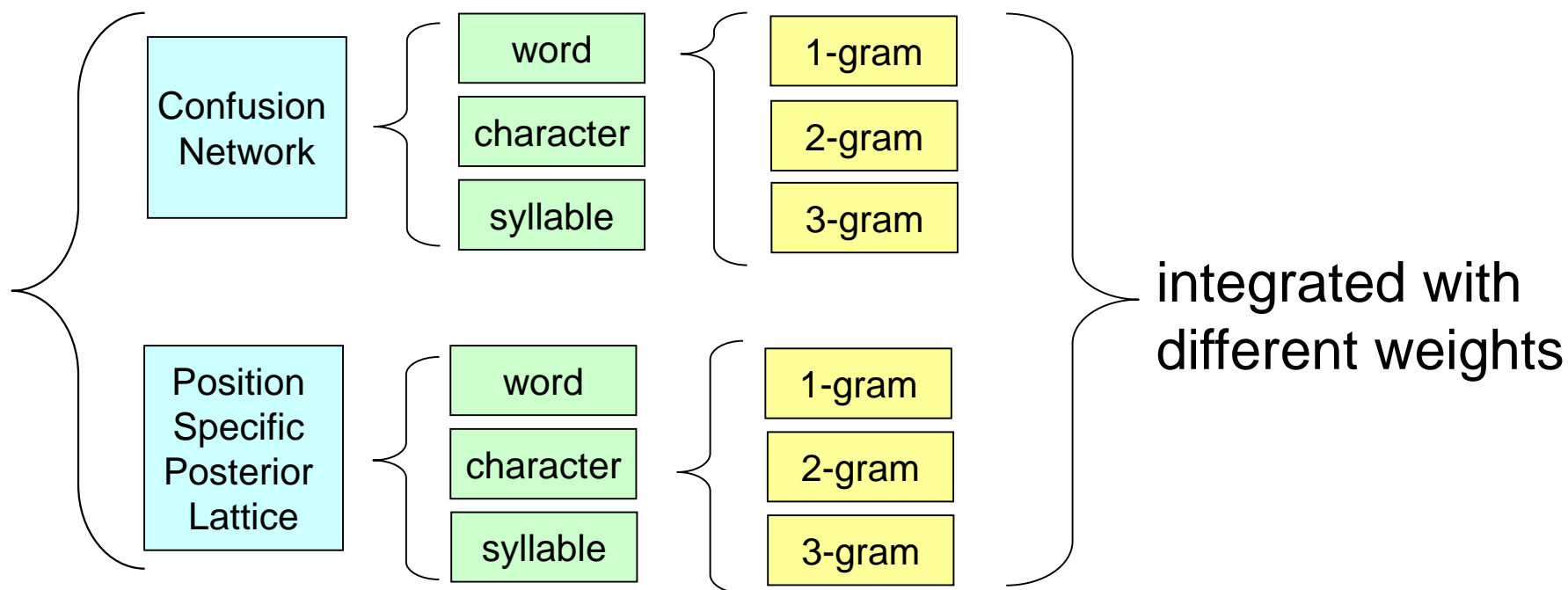
- Linguistically motivated units
 - phonemes, syllables/characters, morphemes, etc.
[Ng, MIT 00][Wallace, Interspeech 07][Chen & Lee, IEEE T. SAP 02]
[Pan & Lee, ASRU 07][Meng, ASRU 07][Meng, Interspeech 08]
[Mertens, ICASSP 09][Itoh, Interspeech 07][Itoh, Interspeech 11]
[Pan & Lee, IEEE T. ASL 10]
- Data-driven units
 - particles, word fragments, phone multigrams, morphs, etc.
[Turunen, SIGIR 07] [Turunen, Interspeech 08]
[Parlak, ICASSP 08][Logan, IEEE T. Multimedia 05]
[Gouvea, Interspeech 10][Gouvea, Interspeech 11][Lee & Lee, ASRU 09]

Integrating Different Clues from Recognition

- Similar to system combination in ASR
- Consistency very often implies accuracy
 - ▣ Integrating the outputs from different recognition systems [Natori, Interspeech 10]
 - ▣ Integrating results based on different subword units [S.-w. Lee, ICASSP 05][Pan & Lee, Interspeech 07][Meng, Interspeech 10][Itoh, Interspeech 11]
 - ▣ Weights of different clues estimated by optimizing some retrieval related criteria [Meng & Lee, ICASSP 09][Chen & Lee, ICASSP 10][Meng, Interspeech 10][Wollmer, ICASSP 09]

Integrating Different Clues from Recognition

- Weights for Integrating 1,2,3-grams for different word/subword units and different indices



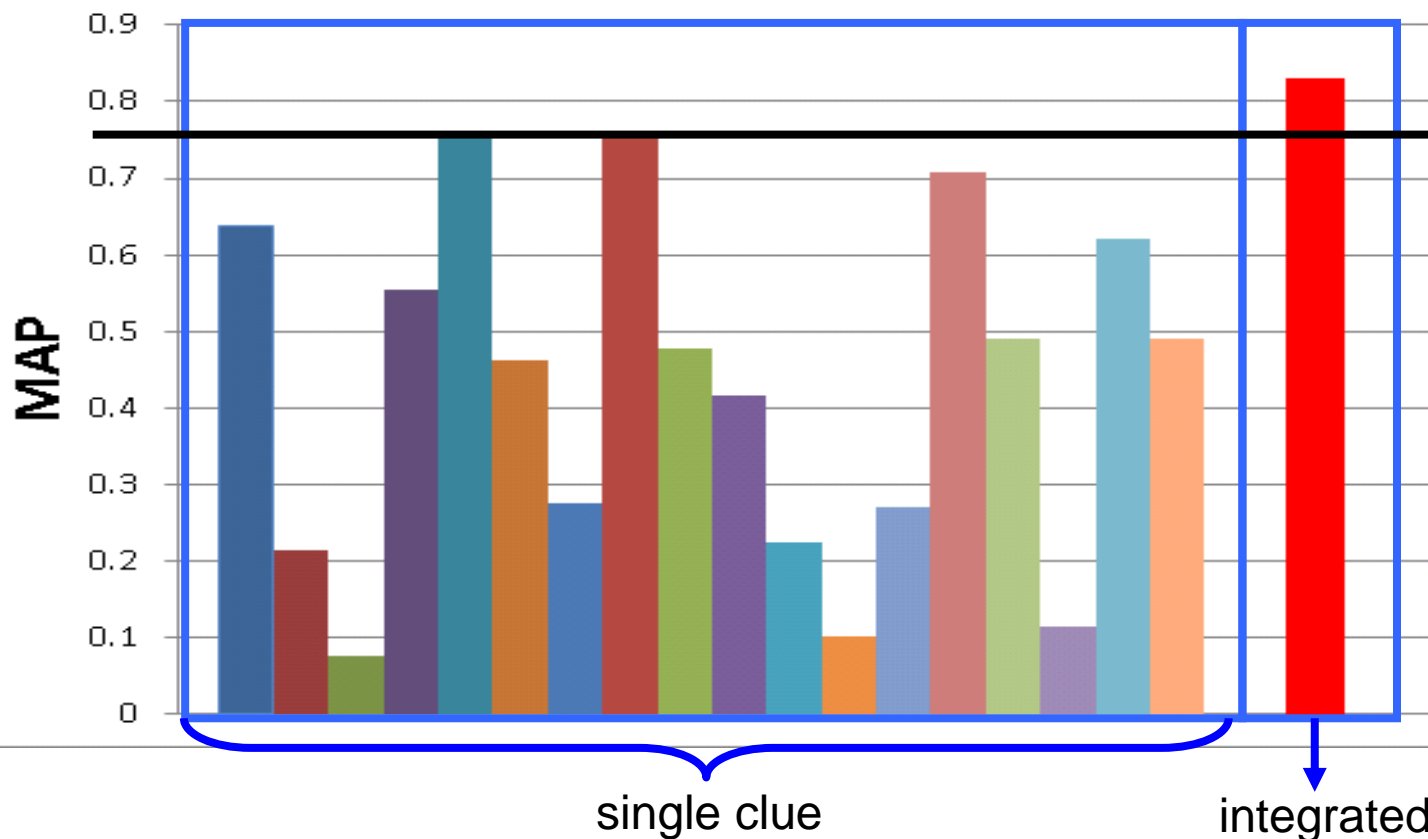
maximizing the lower bound of MAP by SVM-MAP

Training Retrieval Model

Parameters

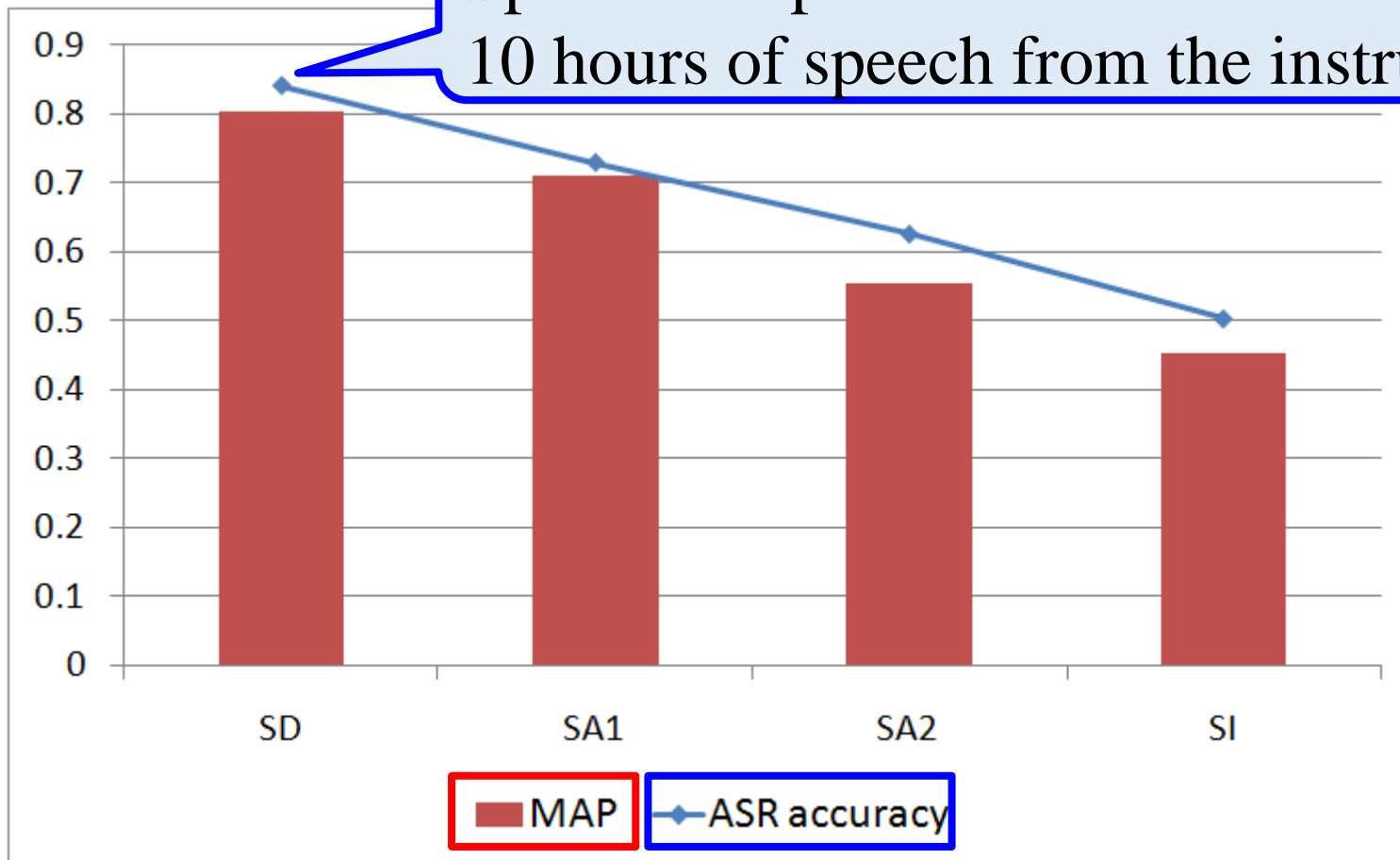
- Integrating different n-grams, word/subword units and indices

[Meng & Lee, ICASSP 09] [Chen & Lee, ICASSP 10]



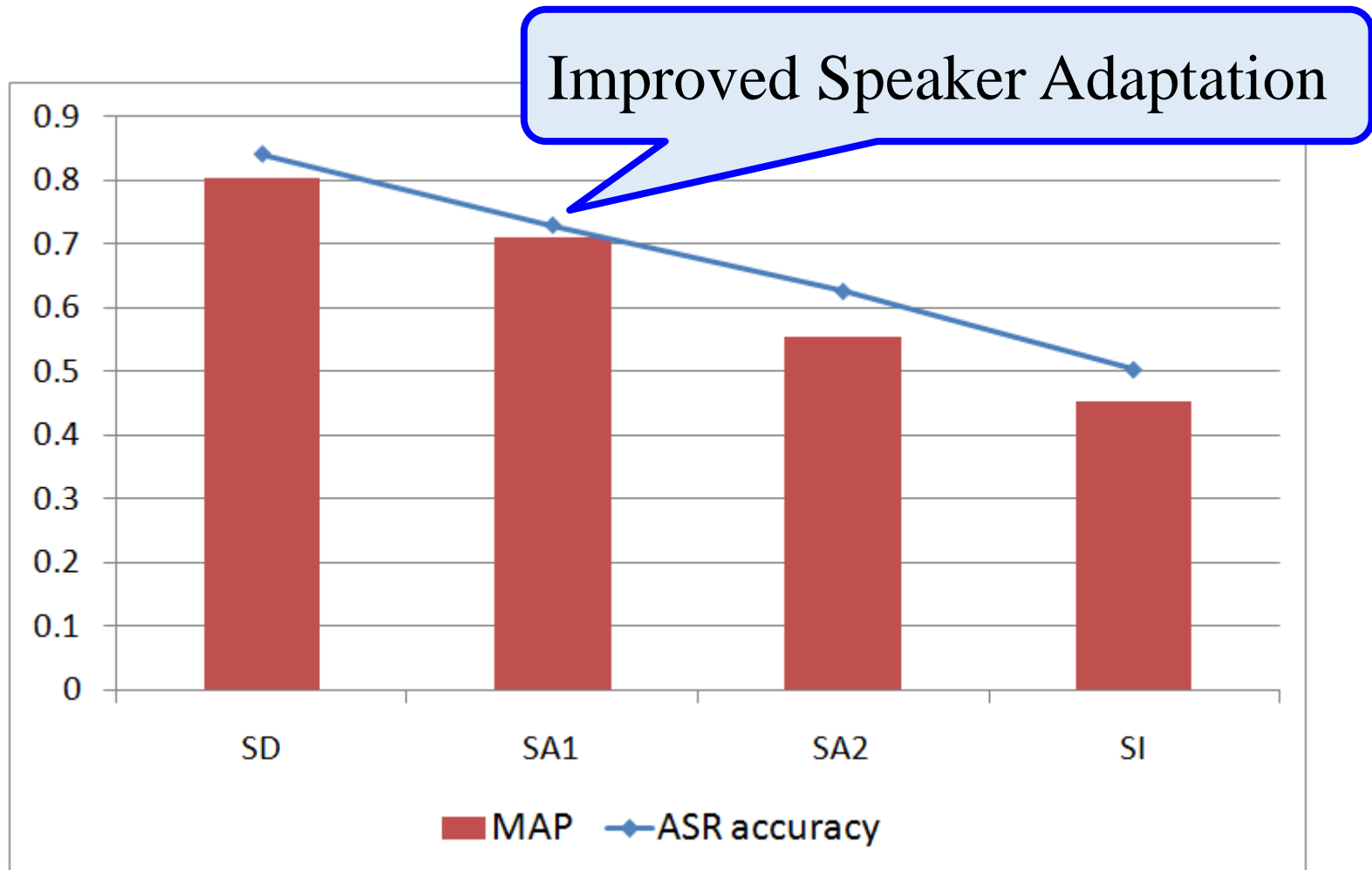
ASR Accuracy v.s. Retrieval Performance

Spoken Term Detection, Lectures



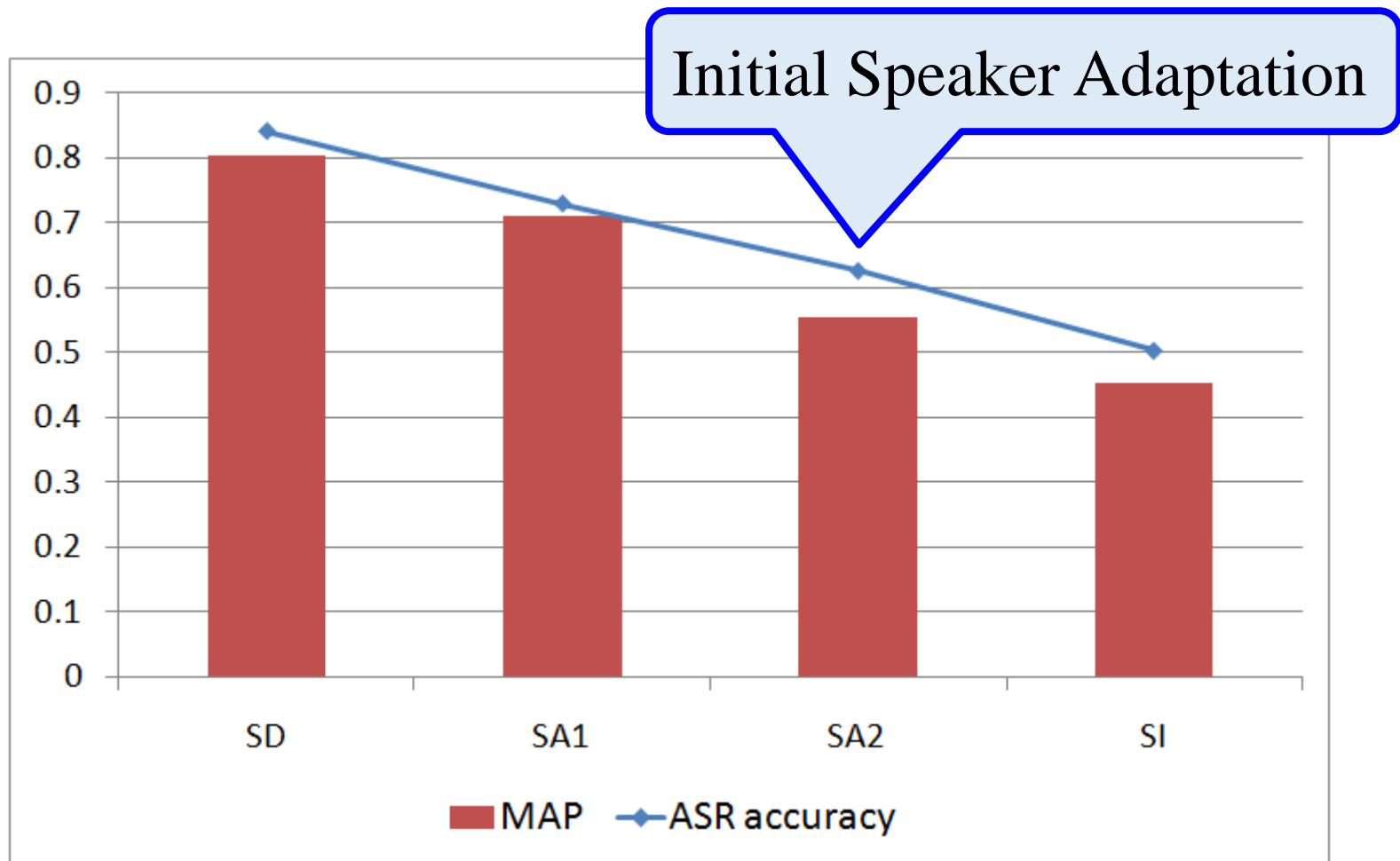
ASR Accuracy v.s. Retrieval Performance

Spoken Term Detection, Lectures



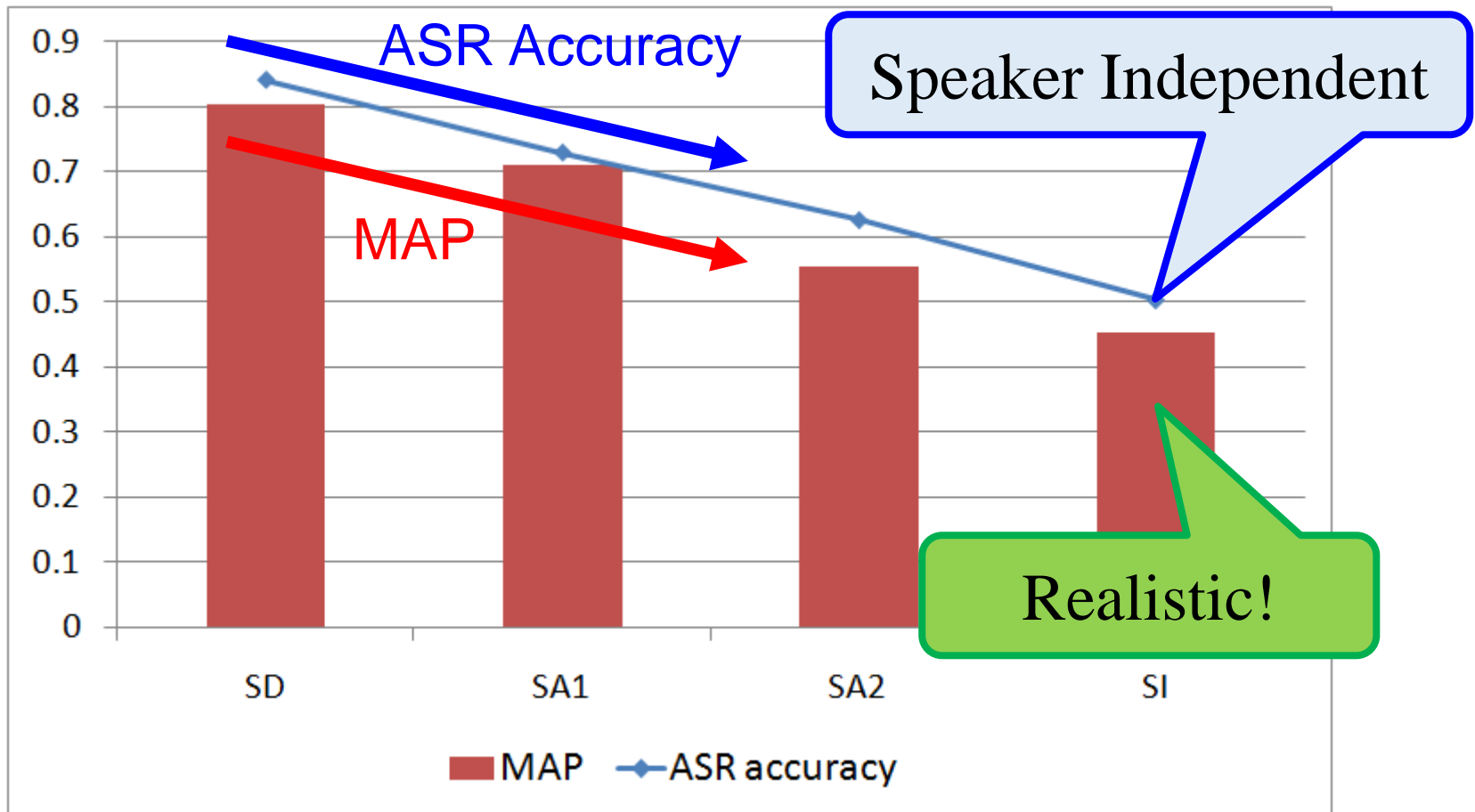
ASR Accuracy v.s. Retrieval Performance

Spoken Term Detection, Lectures



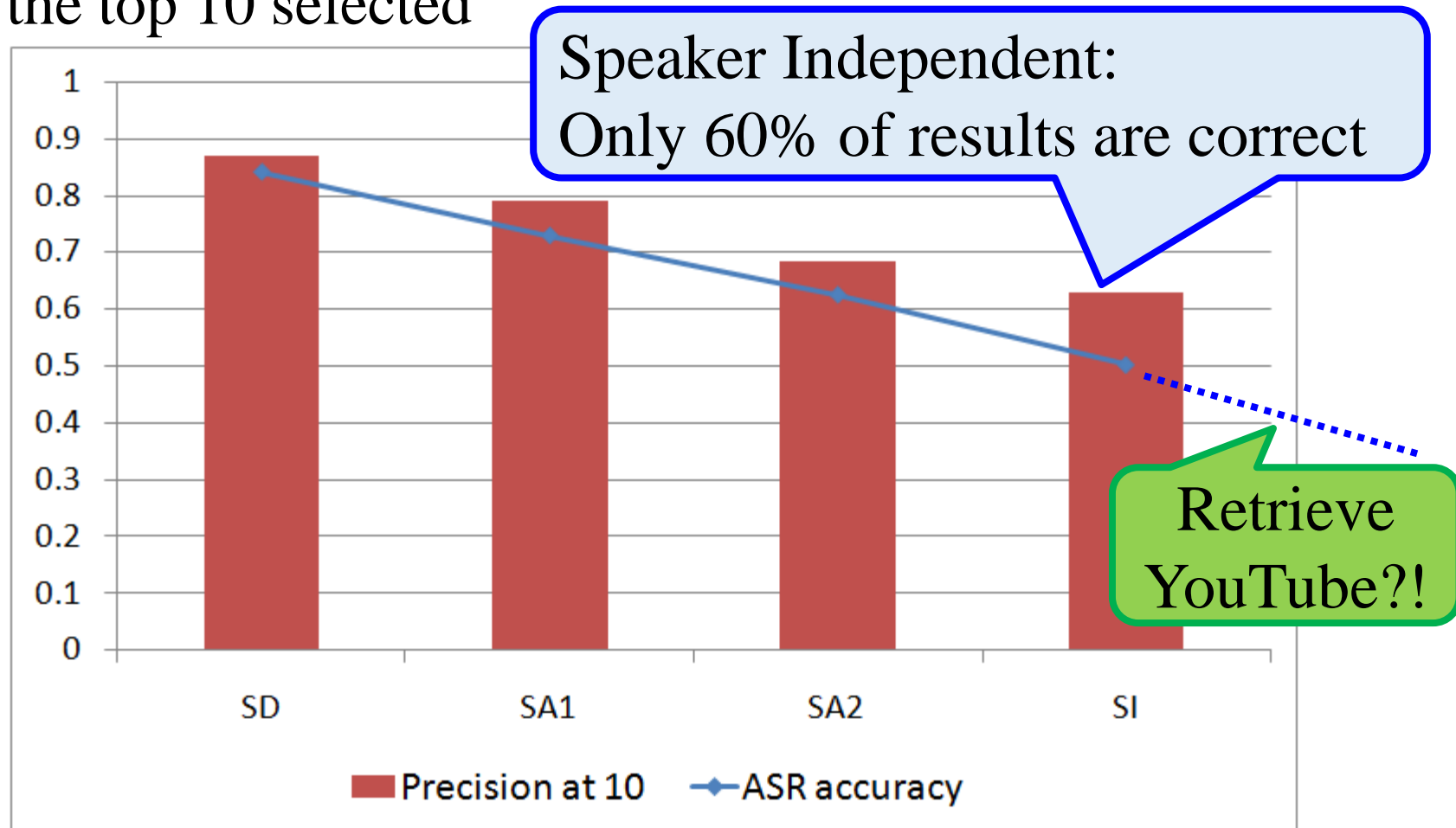
ASR Accuracy v.s. Retrieval Performance

Spoken Term Detection, Lectures



ASR Accuracy v.s. Retrieval Performance

- Precision at 10: Percentage of the correct items among the top 10 selected



Is the problem solved?

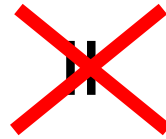
- Did lattices solve the problem?
 - ▣ Need high quality recognition models to produce better lattices and accurately estimate the confidence scores
 - ▣ Spoken content over the Internet is produced in different languages on different domains in different parts of the world under varying acoustic conditions
 - ▣ High quality recognition models for such content doesn't exist yet
- Retrieval performance limited by ASR accuracy

Is the problem solved?

- Desired spoken content retrieval
 - ▣ Less constrained by ASR accuracy
 - ▣ Existing approaches limited by ASR accuracy because of the cascading of speech recognition and text retrieval
- Go beyond the cascading concept

Our point in this tutorial

Spoken Content Retrieval



Speech Recognition

+

Text Retrieval

Core:

Beyond Cascading Speech Recognition and Text Retrieval



New Directions



1. Modified ASR for Retrieval Purposes
2. Incorporating Those Information Lost in ASR
3. No Speech Recognition!
4. Special Semantic Retrieval Techniques for Spoken Content
5. Spoken Content is Difficult to Browse!

Overview Paper

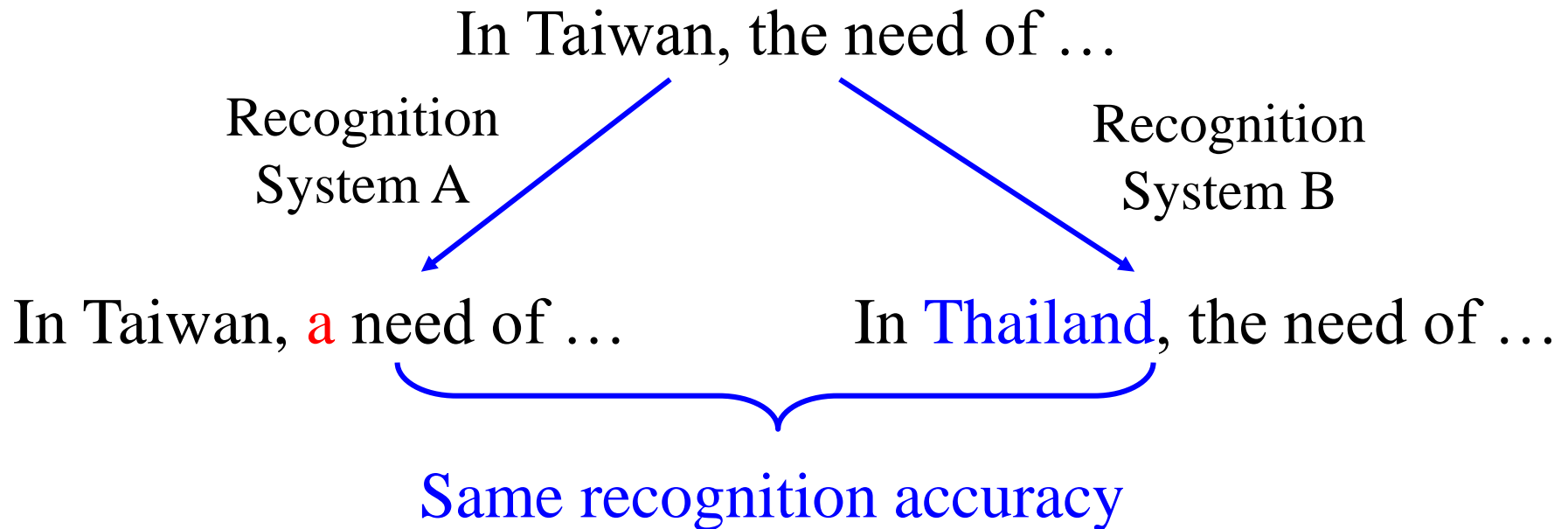
- Lin-shan Lee, James Glass, Hung-yi Lee, Chun-an Chan, "Spoken Content Retrieval —Beyond Cascading Speech Recognition with Text Retrieval," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol.23, no.9, pp.1389-1420, Sept. 2015
- <http://speech.ee.ntu.edu.tw/~tlkagk/paper/Overview.pdf>
- This tutorial includes updated information after this paper is published.

New Direction 1:
Modified ASR
for Retrieval Purposes



Retrieval Performance v.s. Recognition Accuracy

- Intuition: Higher recognition accuracy, better retrieval performance
 - ▣ Not always true!



Retrieval Performance v.s. Recognition Accuracy

- Intuition: Higher recognition accuracy, better retrieval performance
 - ▣ Not always true!

In Taiwan, the need of ...

Recognition
System A

Recognition
System B

In Taiwan, **a** need of ...

Not important
for retrieval

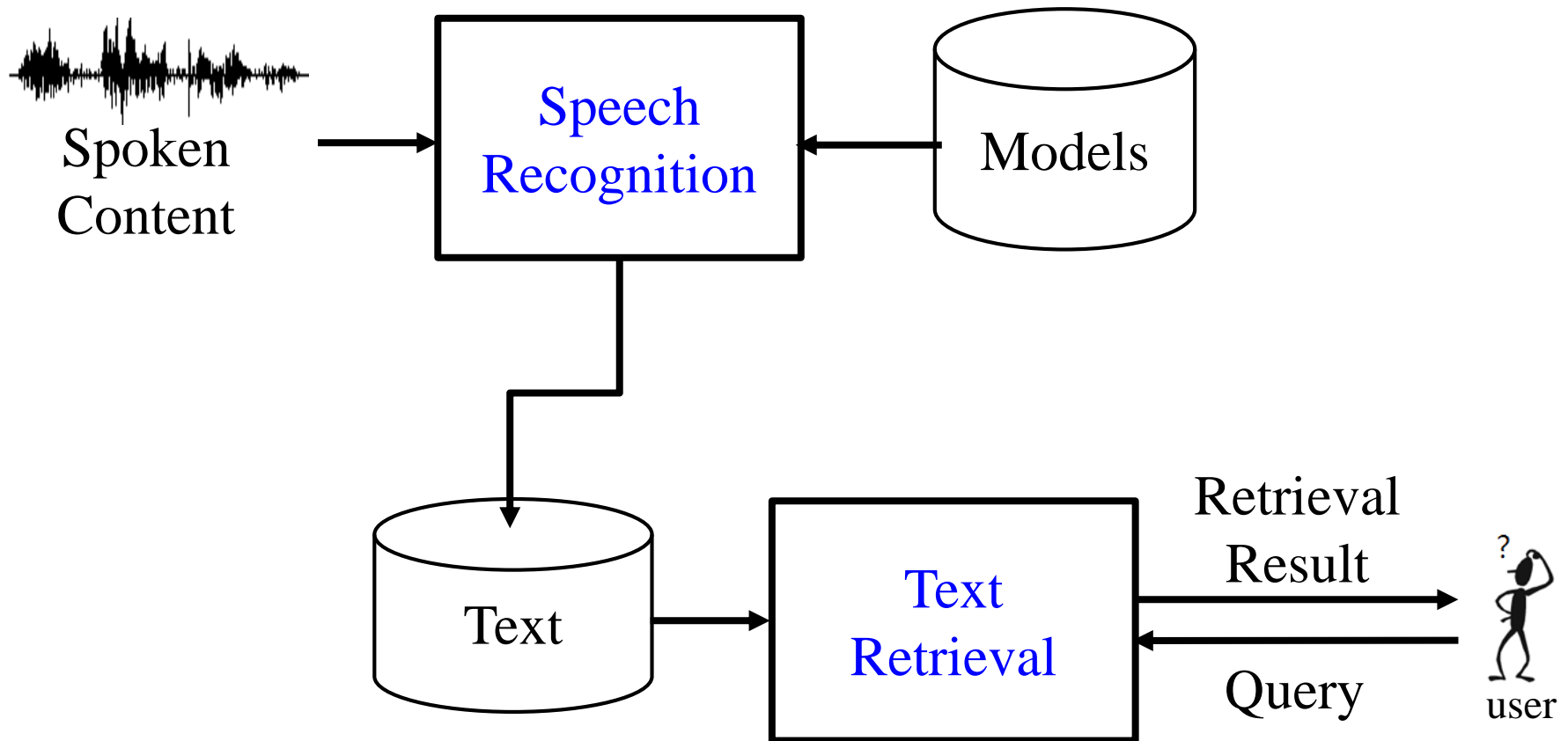
In **Thailand**, the need of ...

Serious problem
for retrieval

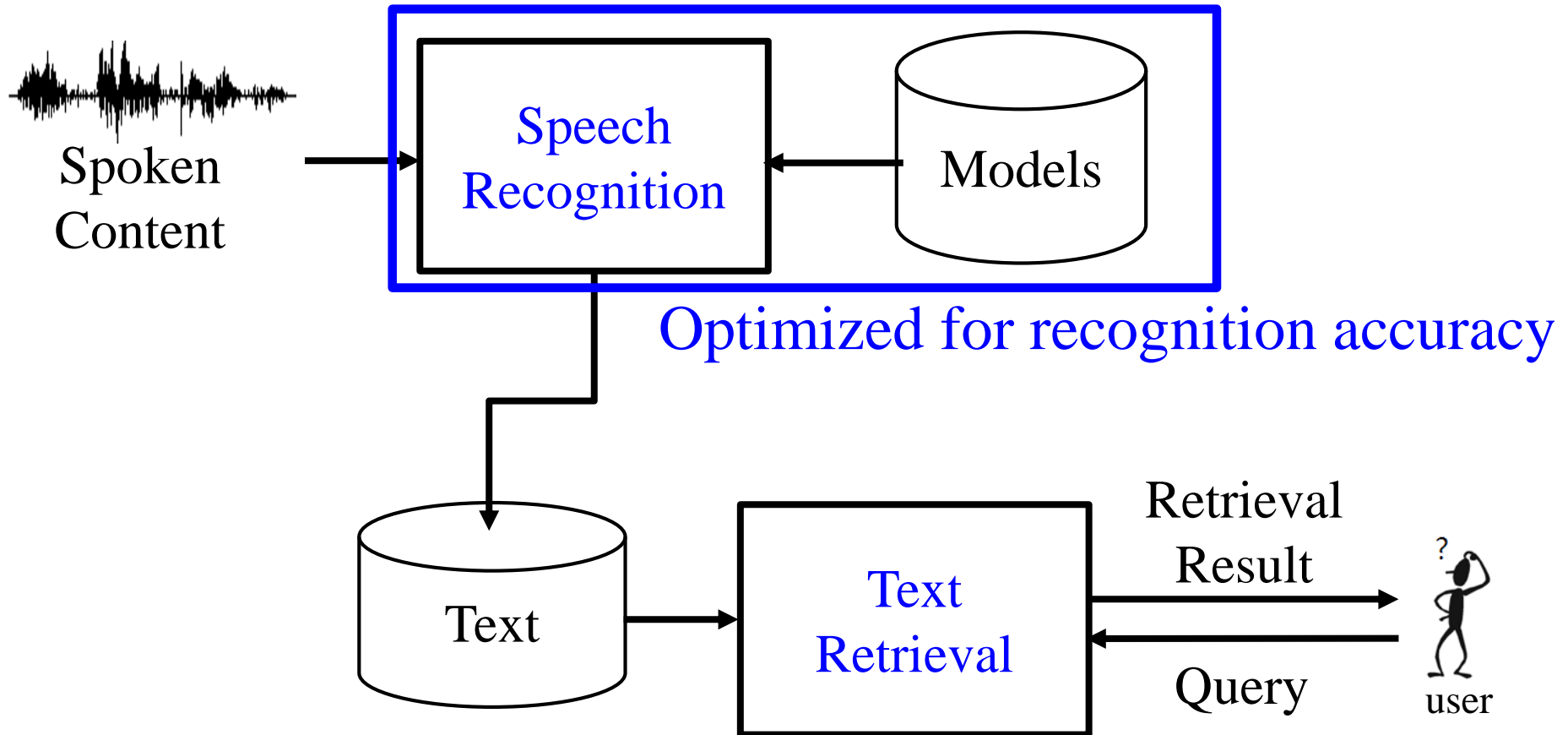
Retrieval Performance v.s. Recognition Accuracy

- Retrieval performance is more correlated to the ASR errors of name entities than normal terms [Garofolo, TREC-7 99][L. van der Werff, SSCS 07]
- Expected error rate defined on lattices is a better predictor of retrieval performance than one-best transcriptions [Olsson, SSCS 07]
 - ▣ lattices used in retrieval
- For retrieval, substitution errors have more influence than insertions and deletions [Johnson, ICASSP 99]
- The language models reducing ASR errors do not always yield better retrieval performance [Cui, ICASSP, 13][Shao, Interspeech, 08][Wallace, SSCS 09]
 - ▣ Query terms usually topic-specific with lower n-gram probabilities

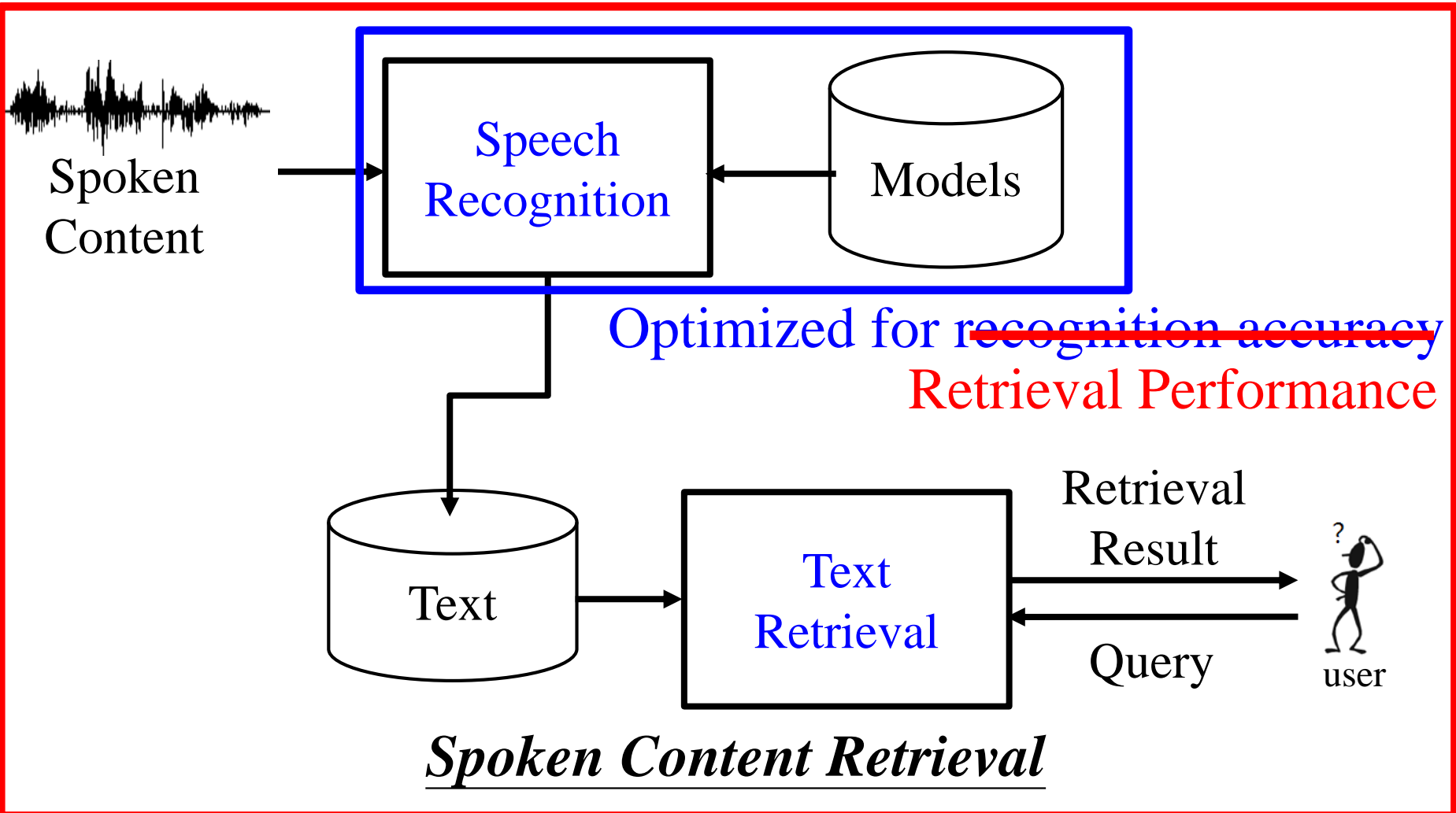
ASR models learned by Optimizing Retrieval Performance



ASR models learned by Optimizing Retrieval Performance



ASR models learned by Optimizing Retrieval Performance



New Direction 1-1:
Modified ASR
for Retrieval Purposes
Acoustic Modeling



Acoustic Modeling

□ Acoustic Model Training

$$\hat{\theta} = \mathit{arg} \max_{\theta} F(\theta)$$

θ : acoustic model parameters
 $F(\theta)$: objective function

The objective function $F(\theta)$ usually defined to optimize ASR accuracy

Design a new objective function for optimizing retrieval performance.

Acoustic Modeling

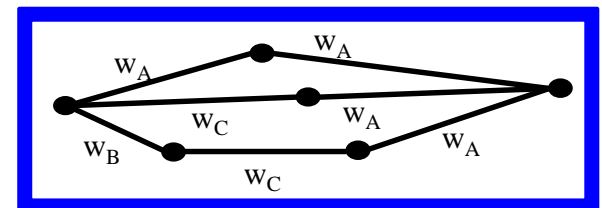
$$\hat{\theta} = \arg \max_{\theta} F(\theta)$$

- Objective Function for optimizing ASR performance

$$F(\theta) = \sum_u \sum_{s_u \in L(u)} A(r_u, s_u) P_{\theta}(s_u | u)$$

Summation over all the utterances u in the training data

➤ $L(u)$: all the word sequence in the lattice of x



lattice of utterance u

Acoustic Modeling

$$\hat{\theta} = \arg \max_{\theta} F(\theta)$$

- Objective Function for optimizing ASR performance

$$F(\theta) = \sum_u \sum_{s_u \in L(u)} A(r_u, s_u) P_{\theta}(s_u|u)$$

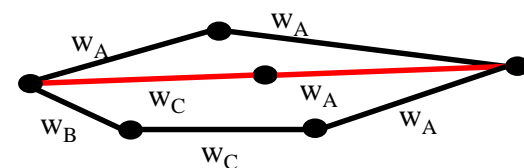
➤ s_u : a word sequence in the lattice of x

➤ $P_{\theta}(s_u|u)$: posterior probability of word sequence s_u given acoustic model θ

➤ $A(r_u, s_u)$: the accuracy of word or phoneme sequence s_u comparing with reference r_u

MCE, MPE, sMBR

θ can be
HMM or DNN



lattice of utterance u

Acoustic Modeling

$$\hat{\theta} = \arg \max_{\theta} F(\theta)$$

- Objective Function for optimizing ~~ASR~~ retrieval performance

$$F(\theta) = \sum_u \sum_{s_u \in L(u)} A(r_u, s_u) P_{\theta}(s_u | u)$$

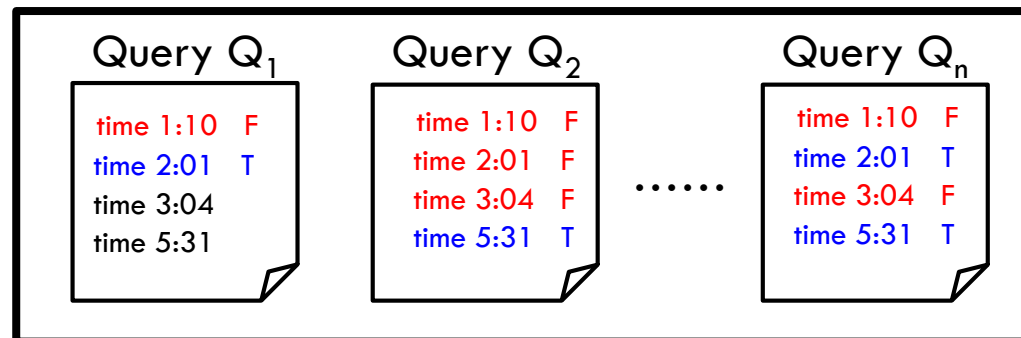
If the possible query terms are known in advance, they can be weighted higher in $A(r_u, s_u)$

W-MCE, [Fu, ASRU 07][Weng, Interspeech 12][Weng, ICASSP, 13]

keyword-boosted sMBR [Chen, Interspeech 14]

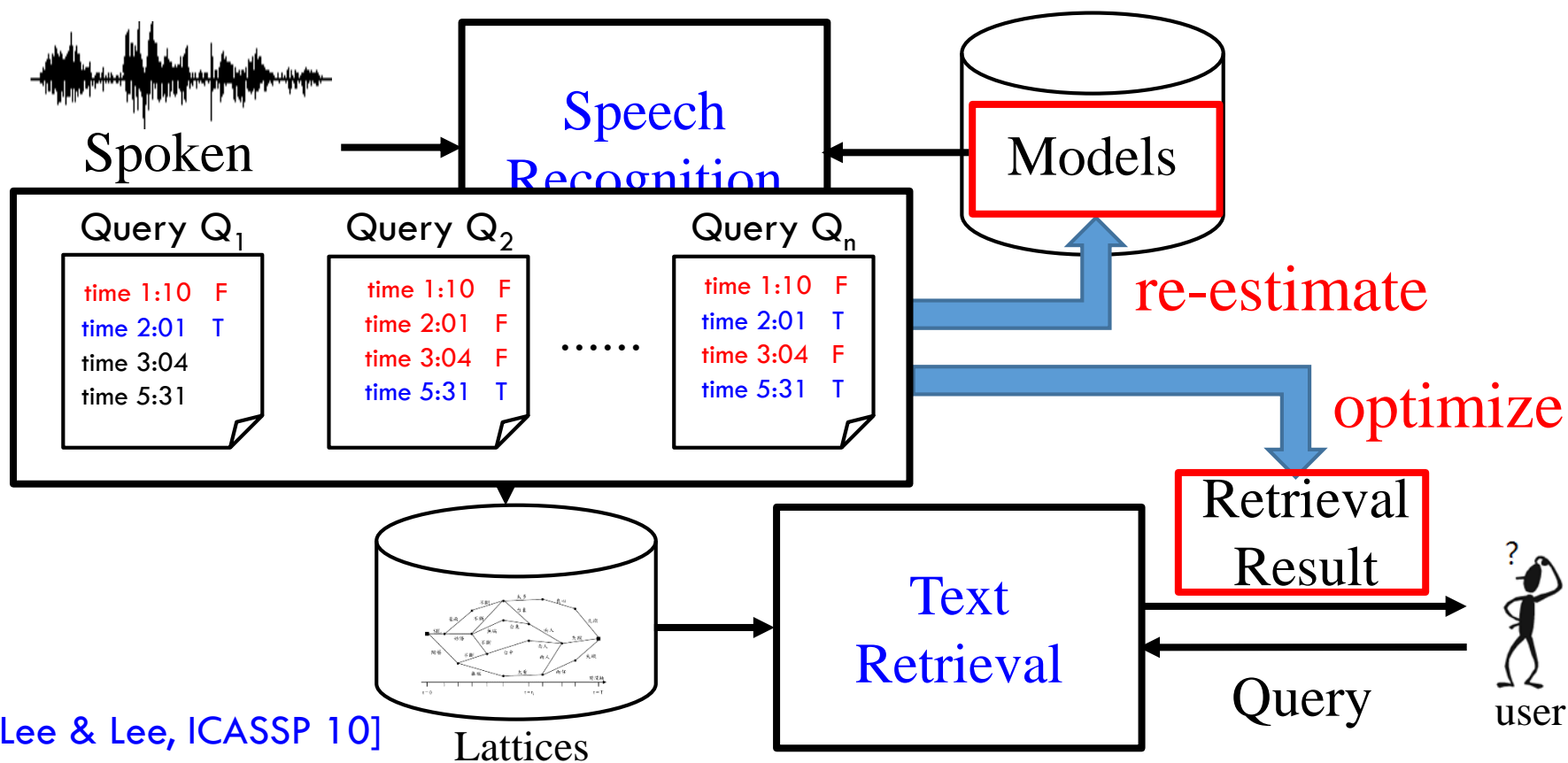
Training Data collected from User

- In most cases, the query terms are not known in advance
- Collect feedback data on-line
 - ▣ Use the information to optimize search engines



- ▣ Feedback can be *implicit*

ASR models learned by Optimizing Retrieval Performance



[Lee & Lee, ICASSP 10]

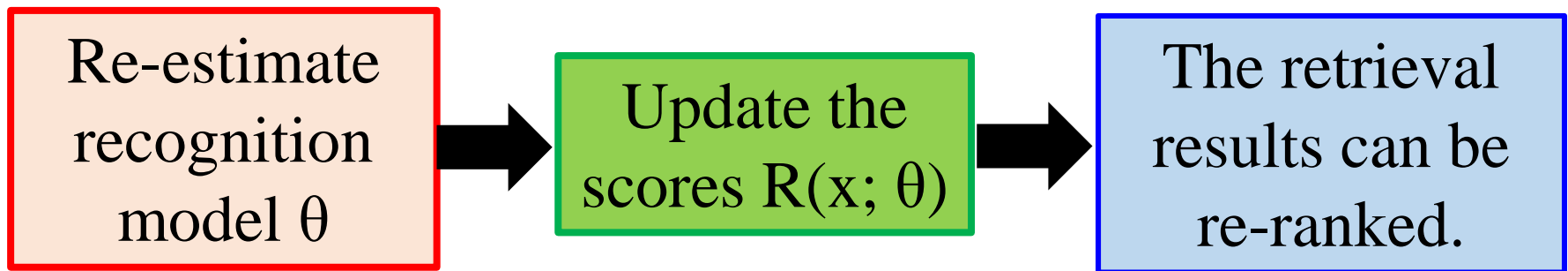
[Lee & Lee, Interspeech 10]

[Lee & Lee, SLT 10]

[Lee & Lee, IEEE T. ASL 12]

Updated Retrieval Process

- Each retrieval result x has a confidence score $R(x)$
- $R(x)$ depends on the recognition model θ
 - $R(x)$ should be $R(x; \theta)$



Considering some
retrieval criterion

Basic Form

$$\hat{\theta} = \arg \max_{\theta} F(\theta)$$

- Basic Form:

$$F(\theta) = \sum_{x_+} R(x_+; \theta) - \sum_{x_-} R(x_-; \theta)$$

x_+ : a positive example

x_- : a negative example

$R(x_+; \theta)$: confidence score of the positive example

$R(x_-; \theta)$: confidence score of the negative example

Basic Form

$$\hat{\theta} = \arg \max_{\theta} F(\theta)$$

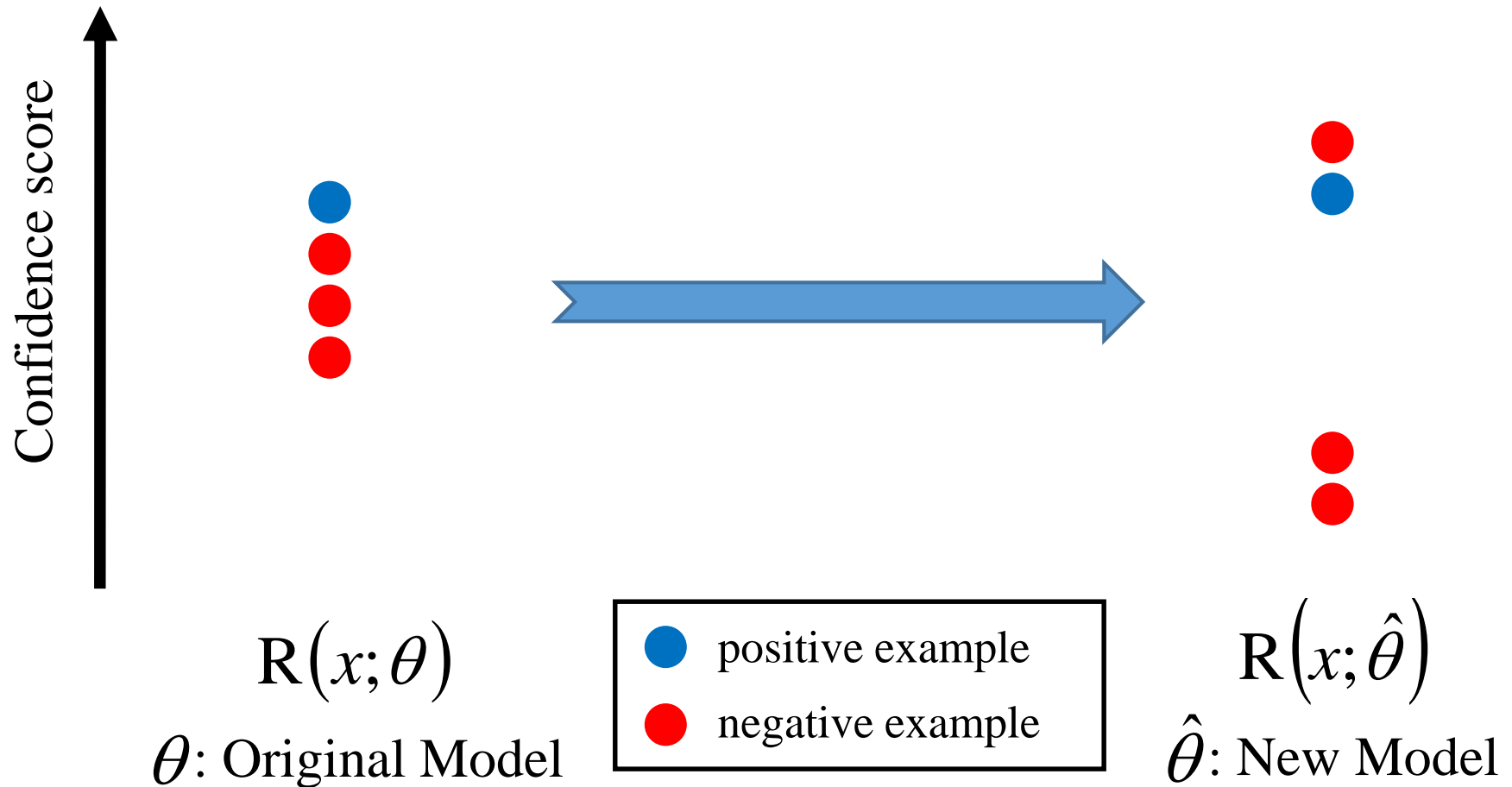
- Basic Form:

$$F(\theta) = \sum_{x_+} R(x_+; \theta) - \sum_{x_-} R(x_-; \theta)$$

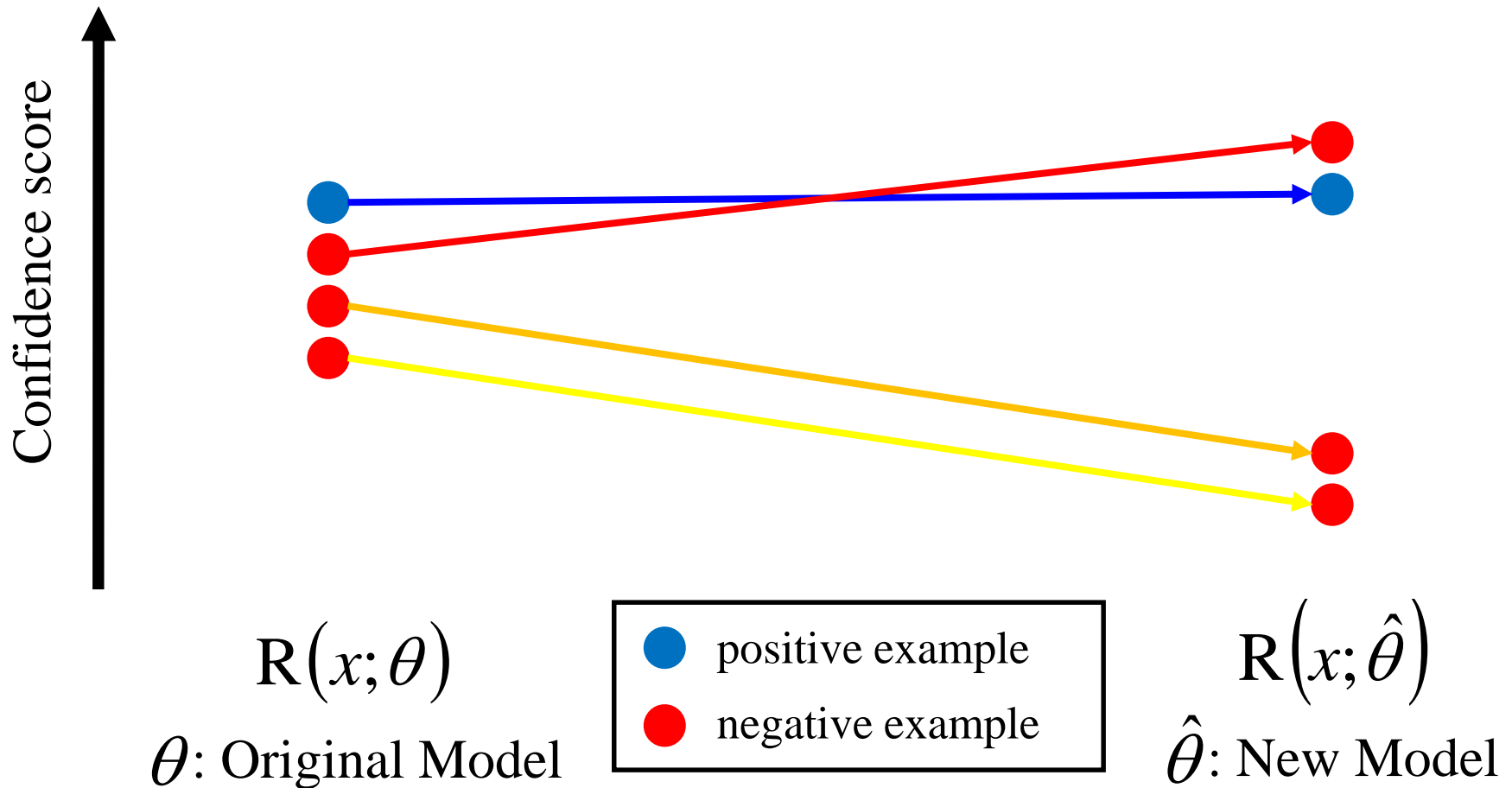
Increase the confidence scores of the positive examples

Decrease the confidence scores of the negative examples

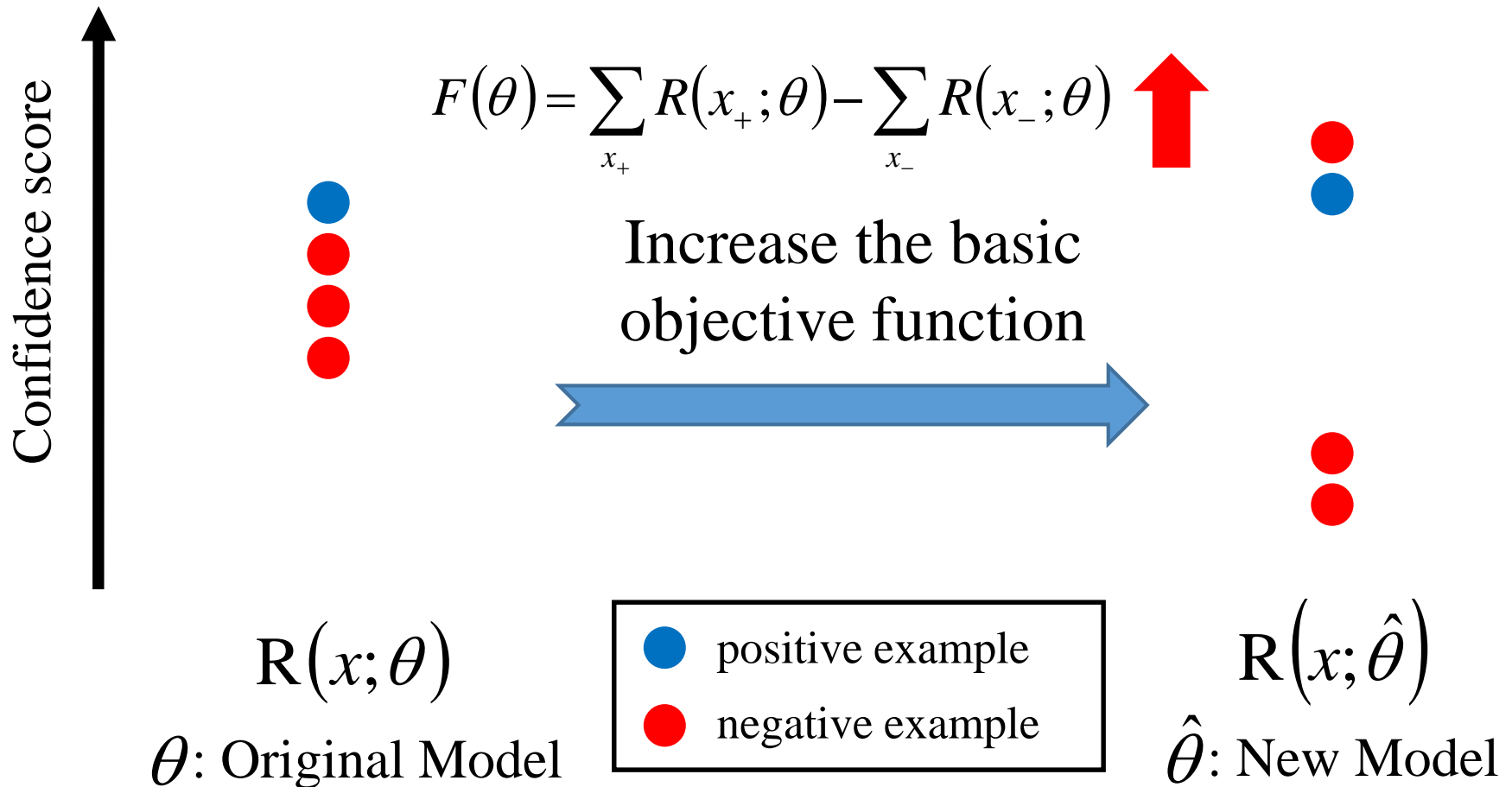
Consider Ranking



Consider Ranking

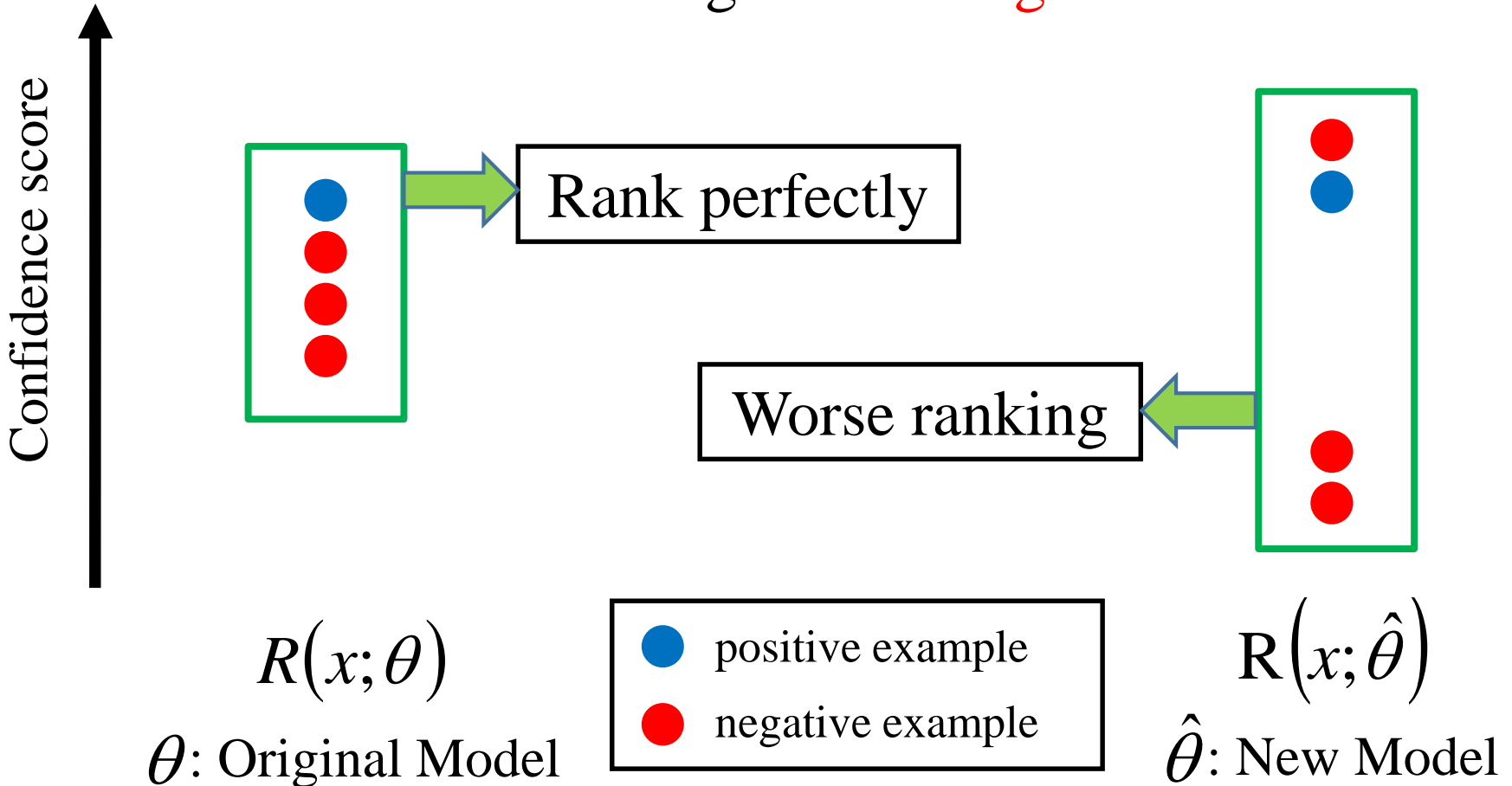


Consider Ranking



Consider Ranking

Considering the **ranking order**



Consider Ranking

$$F(\theta) = \sum_{x_+, x_-} \delta(x_+, x_-)$$

$$\delta(x_+, x_-) = \begin{cases} 1 & \mathbf{R}(x_+; \theta) > \mathbf{R}(x_-; \theta) \\ 0 & \textit{otherwise} \end{cases}$$

- If the confidence score for a positive example exceed that for a negative example
 - the objective function adds 1.

Consider Ranking

$$F(\theta) = \sum_{x_+, x_-} \delta(x_+, x_-)$$

$$\delta(x_+, x_-) = \begin{cases} 1 & \mathbf{R}(x_+; \theta) > \mathbf{R}(x_-; \theta) \\ 0 & \textit{otherwise} \end{cases}$$

- $\delta(x_+, x_-)$ approximated by a sigmoid function during optimization.

Little feedback data?

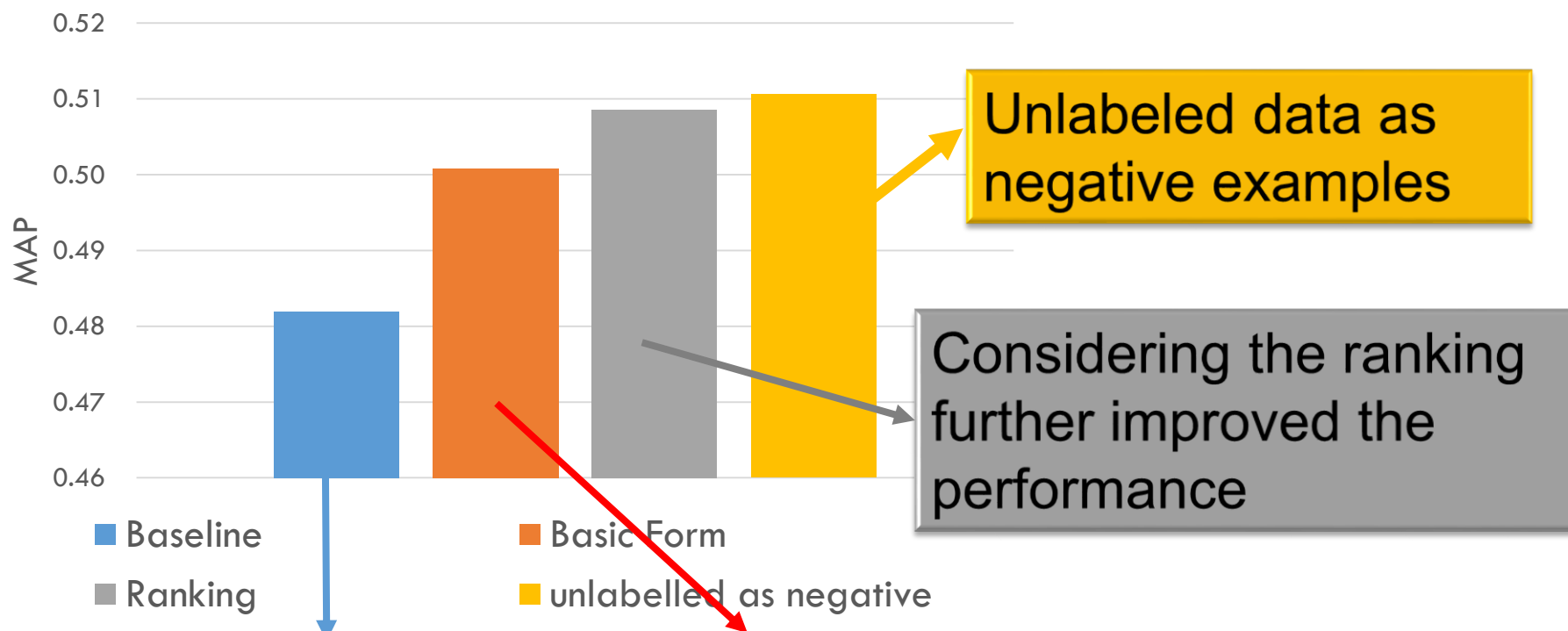


The unlabeled examples as negative examples

Acoustic Models - Experiments

- Lecture recording (80 queries, each has 5 clicks)

[Lee & Lee, IEEE T. ASL 12]



Baseline: before model re-estimation

Re-estimating acoustic models improved the performance (basic form)

New Direction 1-2:
Modified ASR
for Retrieval Purposes
Besides Acoustic Modeling



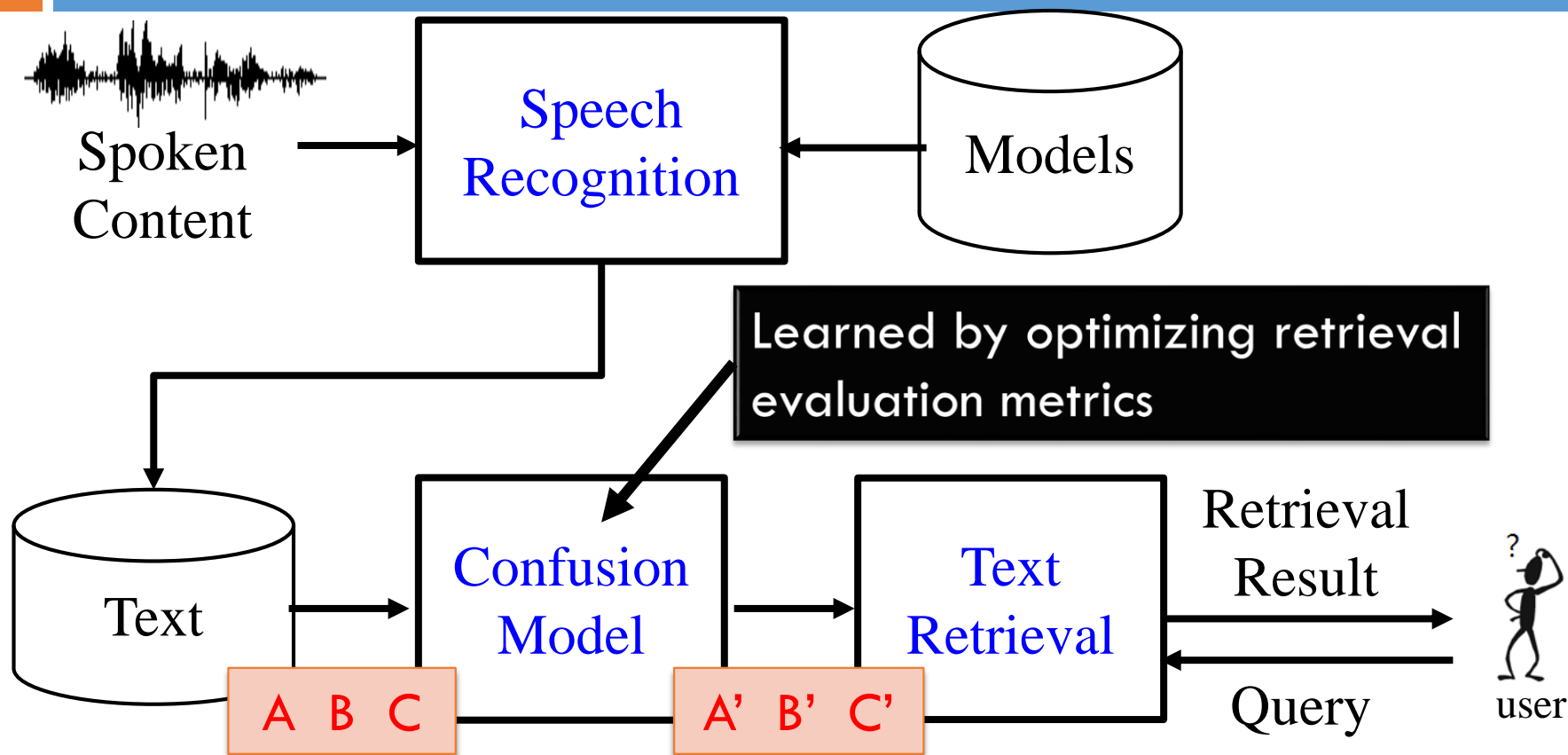
Language Modeling

- The query terms are usually very specific. Their probabilities are underestimated.
- Boosting the probabilities of n-grams including query terms
 - ▣ By repeating the sentences including the query terms in training corpora
 - ▣ Helpful in DARPA's RATS program [Mandal, Interspeech 13] and NIST OpenKWS13 evaluation [Chen, ISCSLP 14]
- NN-based LM: Modifying training criterion, so the key terms are weighted more during training
 - ▣ Helpful in NIST OpenKWS13 evaluation [Gandhe, ICASSP 14]

Decoding

- Give different words different pruning thresholds during decoding
 - ▣ The keywords given lower pruning thresholds than normal terms
 - ▣ Called white listing [Zhang, Interspeech 12] or keyword-aware pruning [Mandal, Interspeech 13]
- OOV words never correctly recognized
 - ▣ Two stage approach [Shao, Interspeech 08]
 - Identify the lattices probably containing OOV (by subword-based approach)
 - Insert the word arcs of OOV words into lattices and rescore

Confusion Models



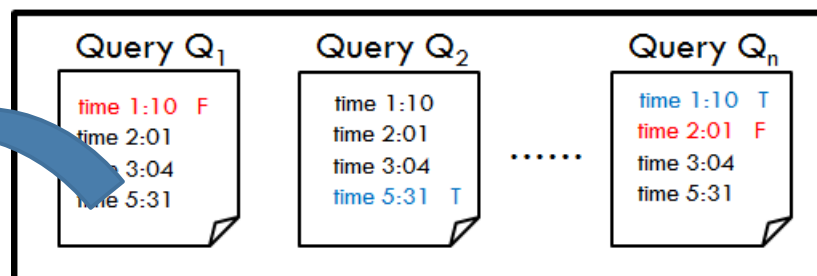
The ASR produces systematic errors, so it is possible to learn a confusion model to offer better retrieval results

[Karanasou, Interspeech 12][Wallace, ICASSP 10]

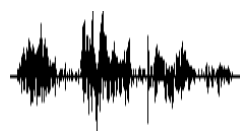
Jointly Optimizing Speech Recognition and Retrieval Modules

Sounds crazy?

Structured SVM with
Hidden variables
[R. Prabhavalkar, ICASSP, 2013,
MLSLP 2012]



A spoken segment



query

Complex
Model

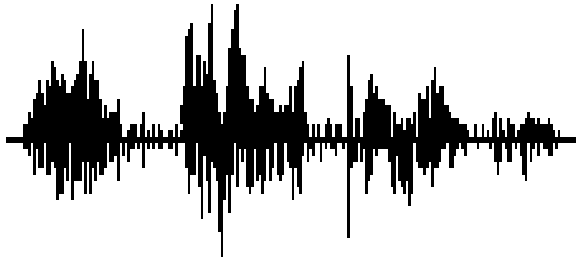
Yes, the segment
contains the query.

No,

End-to-end model performing speech recognition and retrieval jointly (learned jointly) in one step

Much information lost during ASR

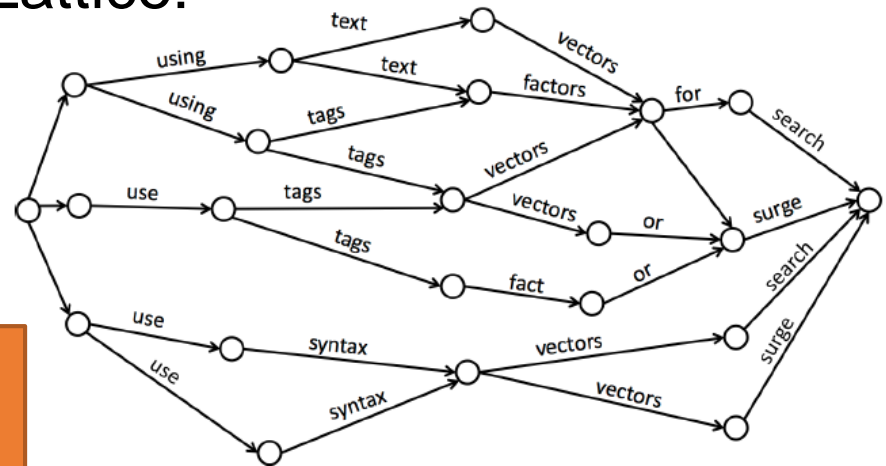
Spoken Content



ASR

Transcriptions:
using syntax vectors surge

Lattice:



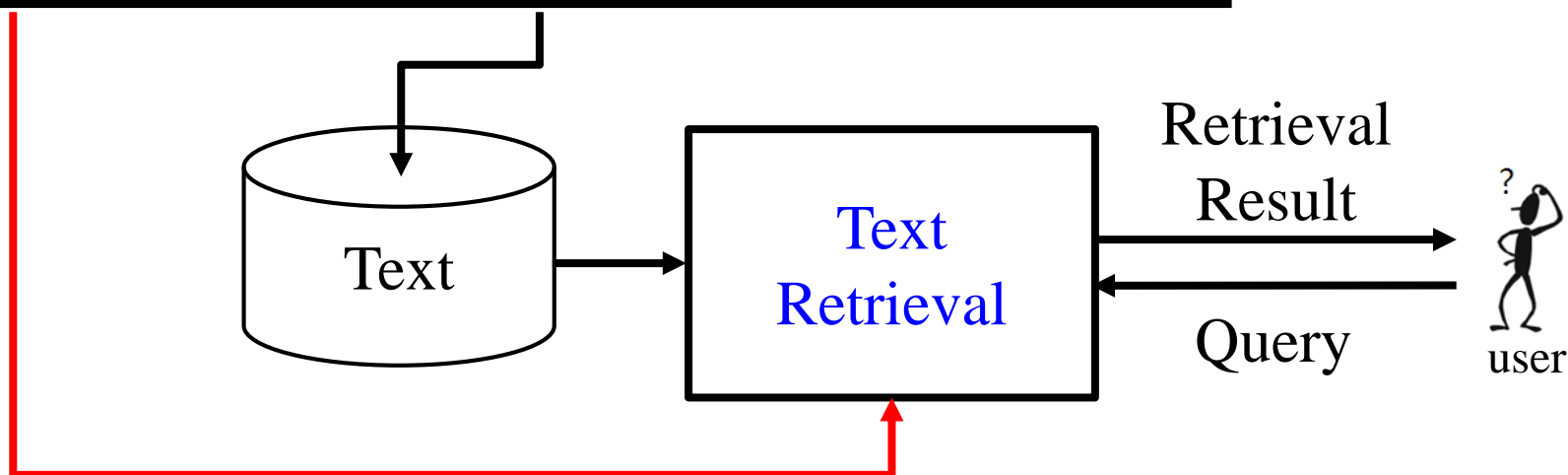
Much information lost during ASR

New Direction 2:
Incorporating
Those Information Lost in ASR



Information beyond Speech Recognition Output

Black Box



Incorporating information lost in
ASR to help retrieval

New Direction 2-1:

Incorporating

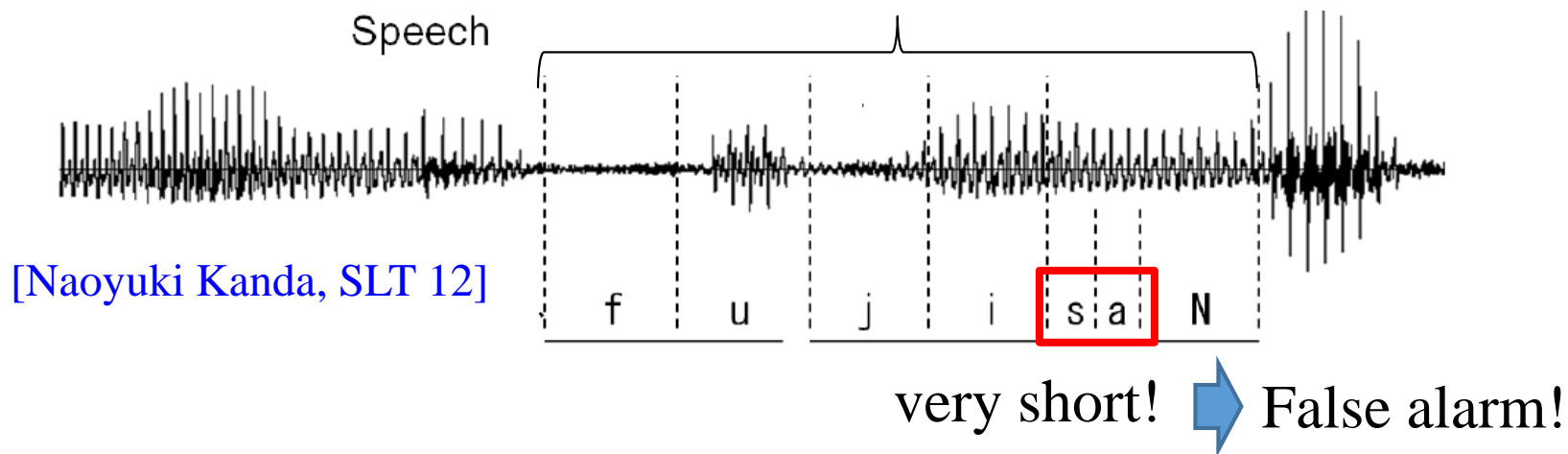
Those Information Lost in ASR

What kind of information can be helpful?

Information beyond Speech Recognition Output

- Phoneme or syllable duration [Wollmer, ICASSP 09][Naoyuki Kanda, SLT 12][Tepei Ohno, SLT 12]

Query is Japanese word “fu-ji-sa-N”



- Pitch & Energy [Tejedor, Interspeech 10]
- Landmark and attribute detection with prosodic cues includes can reduce the false alarm [Ma, Interspeech 2007]

Query-specific Information

- **"Jack of all trades, master of none"**

Speech Recognition



Correctly recognized
all the words

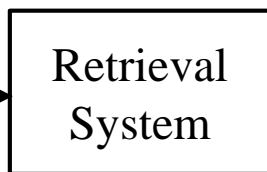
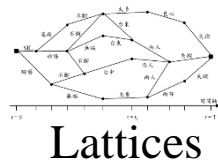
Spoken Term
Detection

Query-specific
detector



higher detector accuracy
on specific query

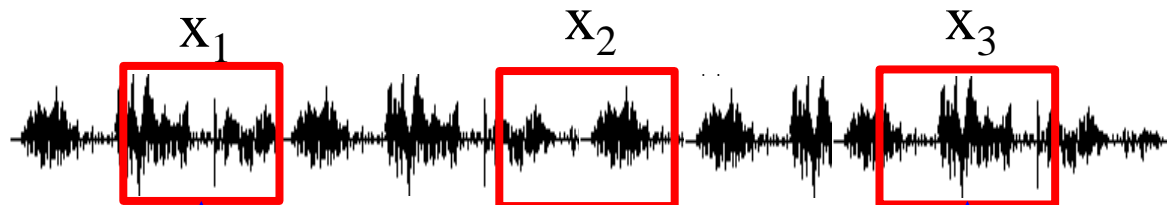
Query-specific Detector



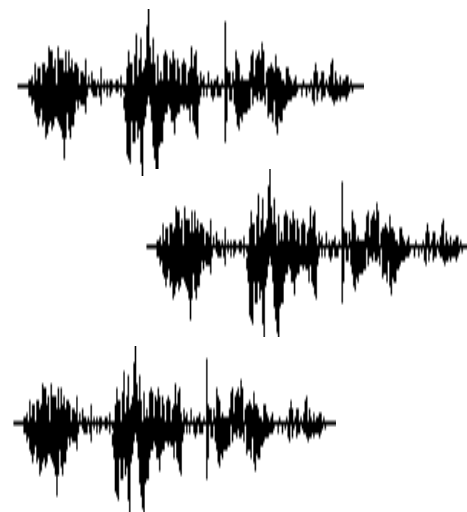
Query Q



First-pass Retrieval Result



Examples of Q



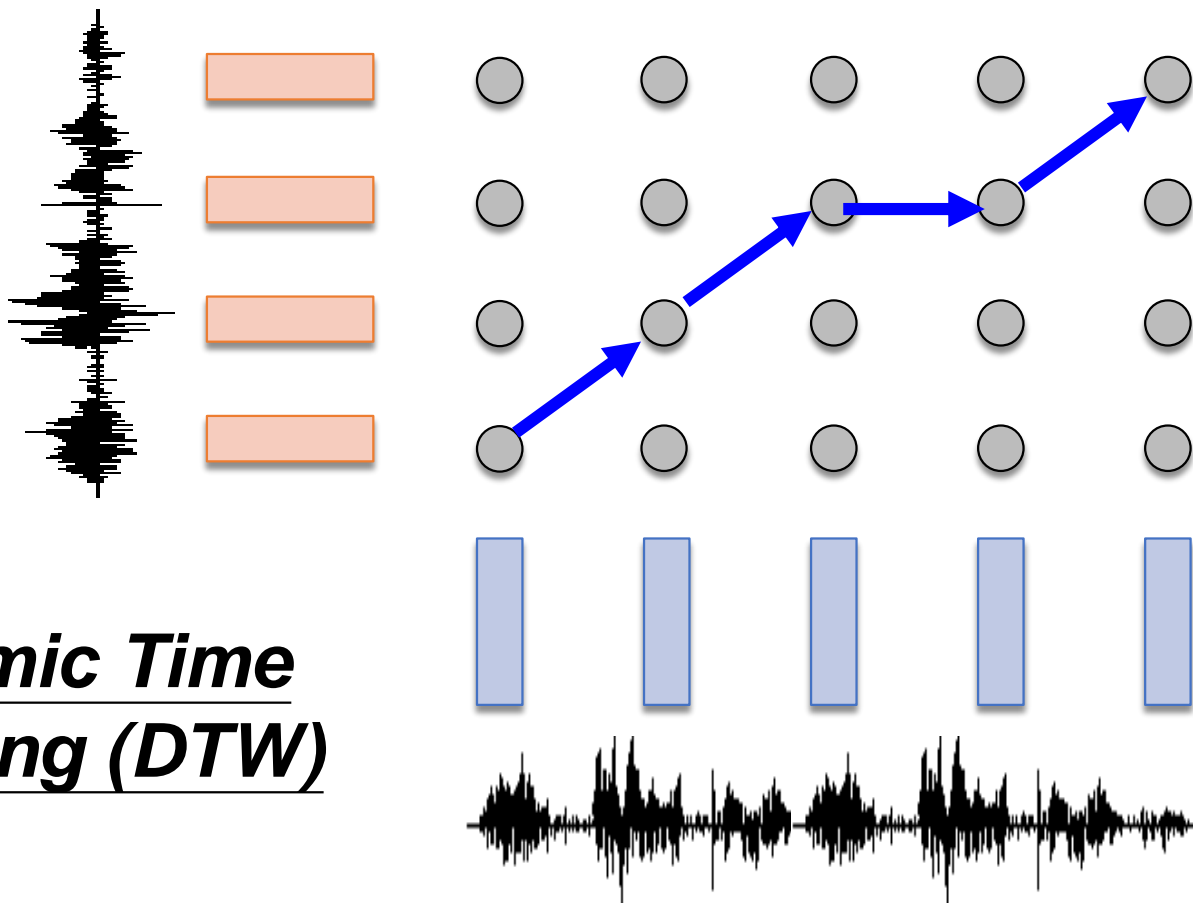
Compute Similarity

Exemplar-based approach also used in speech recognition

[Demuyck, ICASSP 2011][Heigold, ICASSP 2012][Nancy Chen, ICASSP 2016]

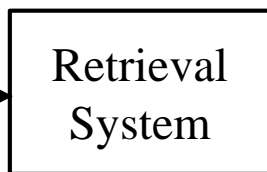
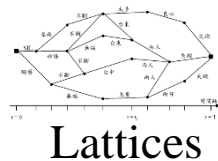
Similarities

between Audio Segments



Dynamic Time
Warping (DTW)

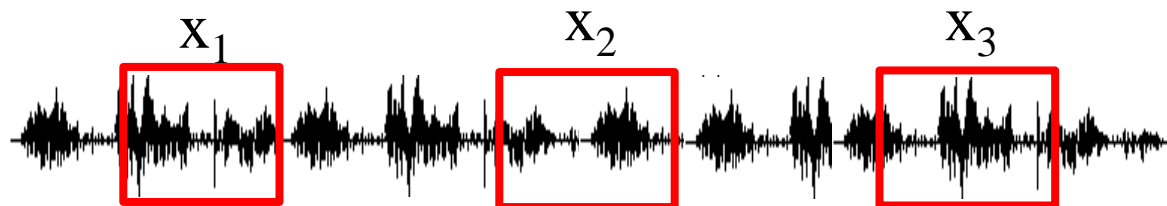
Query-specific Detector



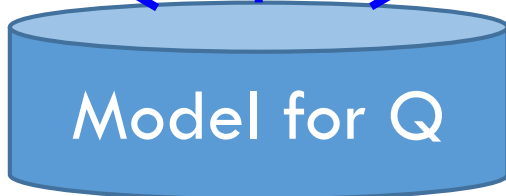
Query Q



First-pass Retrieval Result



Evaluate confidence

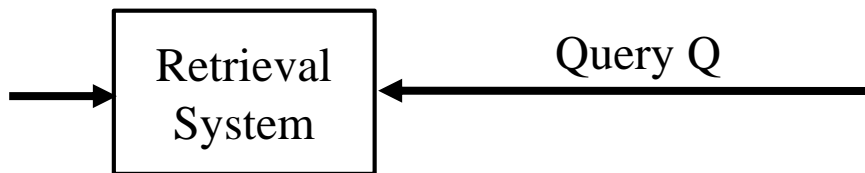
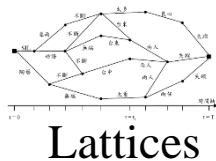


Examples of Q



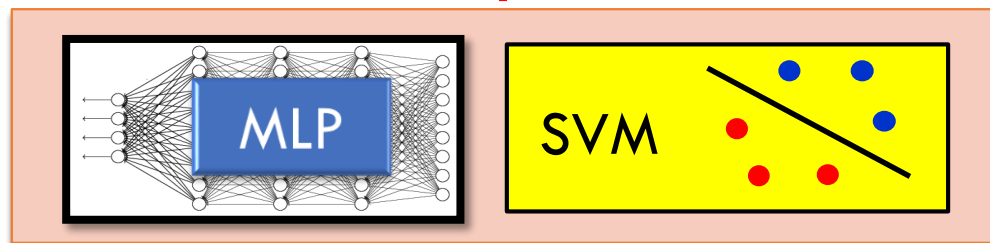
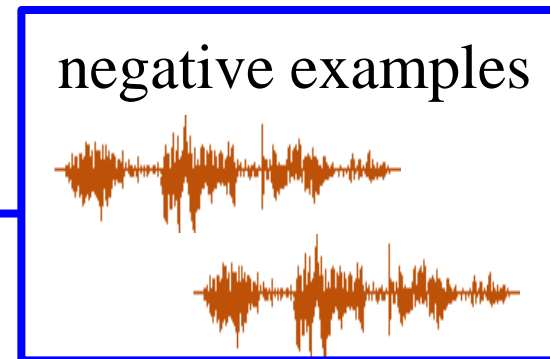
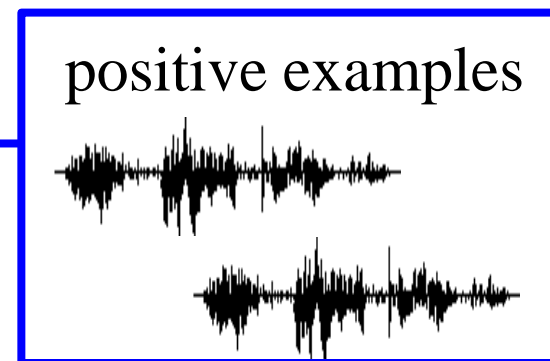
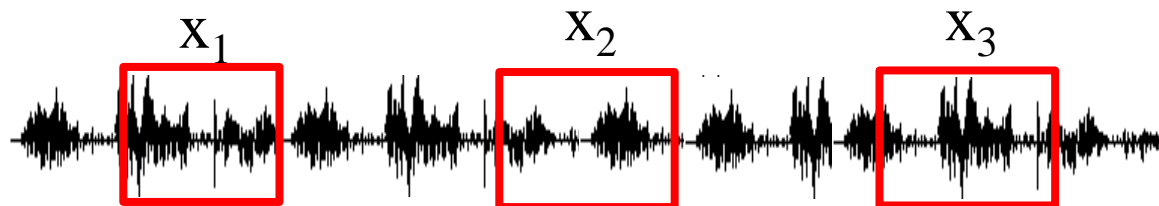
Learn a model

Query-specific Detector



[Tu & Lee, ASRU 11]
[I.-F. Chen, Interspeech 13]

First-pass Retrieval Result

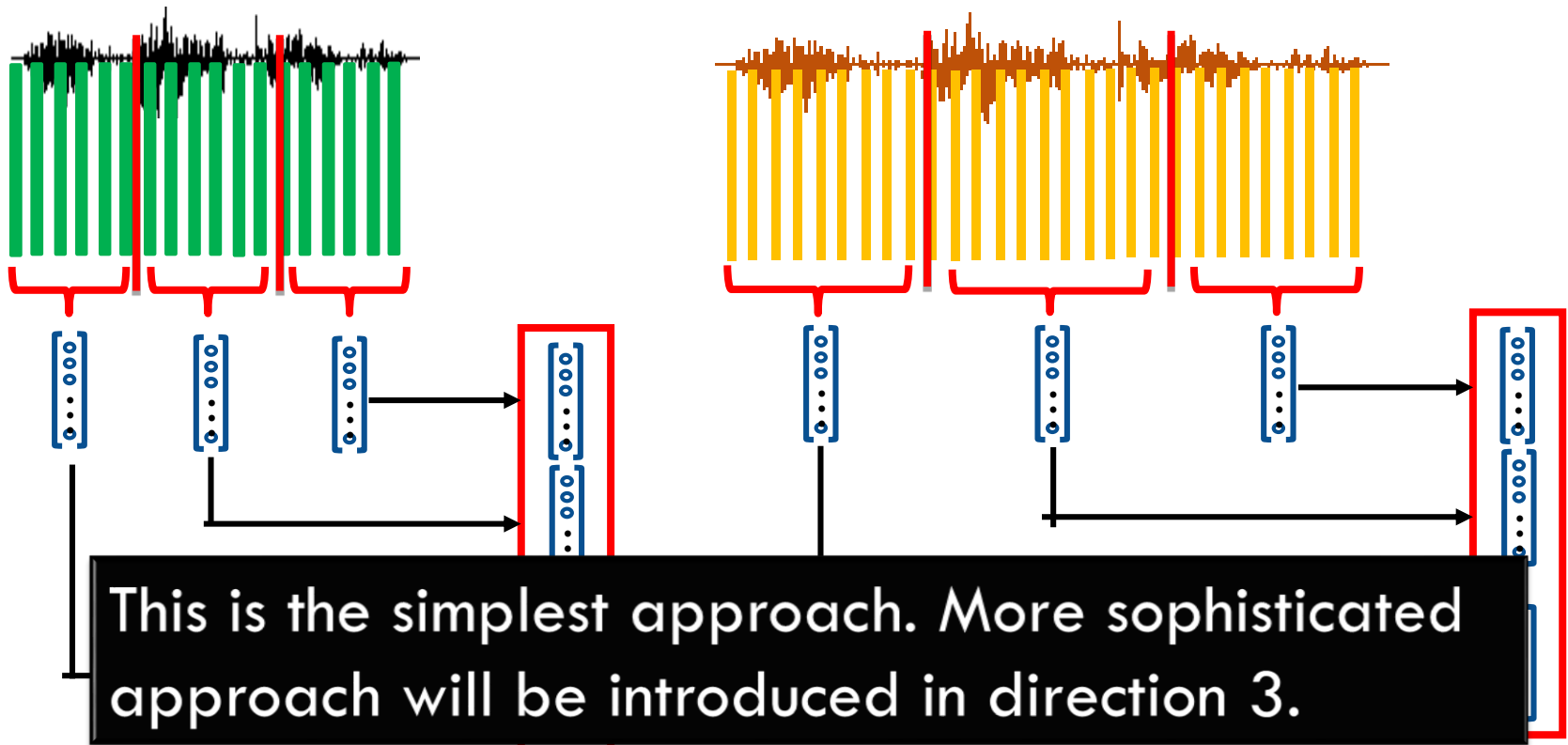


Learn a discriminative model

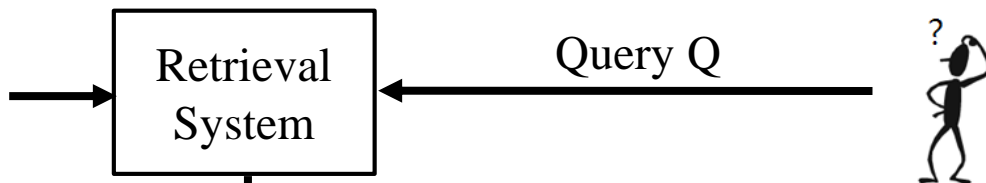
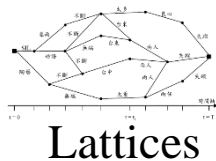
Query-specific Detector

[Tu & Lee, ASRU 11]
[I.-F. Chen, Interspeech 13]

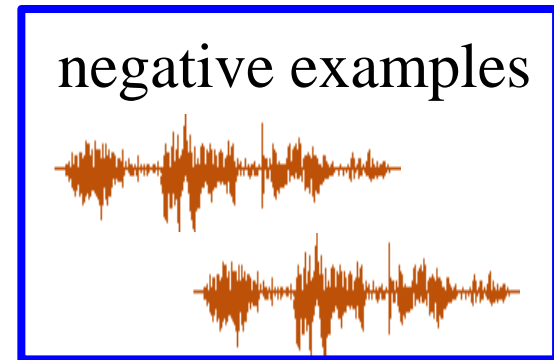
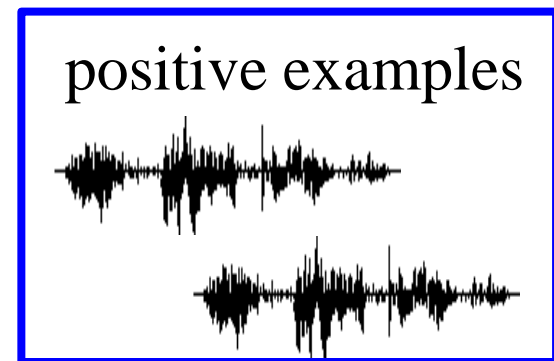
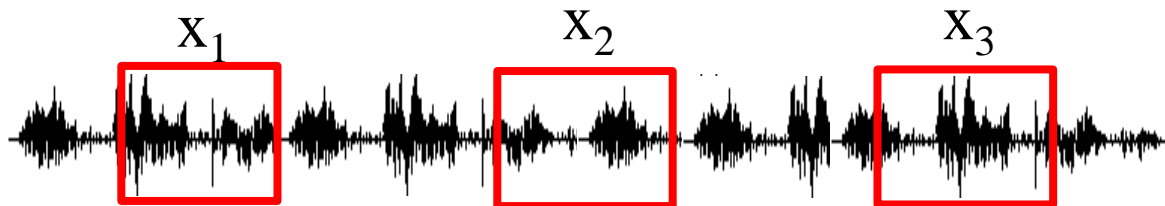
- The input of SVM or MLP has to be a fixed-length vector
- Representing an audio segment with different length into a fixed-length vector



Query-specific Detector



First-pass Retrieval Result



➤ Is it realistic to have those examples?

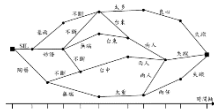
➡ Data collected from users (direction 1)

➡ Pseudo-relevance Feedback (PRF)

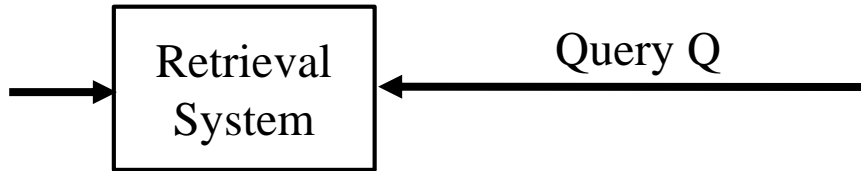
New Direction 2-2:
Incorporating
Those Information Lost in ASR
Pseudo Relevance Feedback



Pseudo Relevance Feedback (PRF)

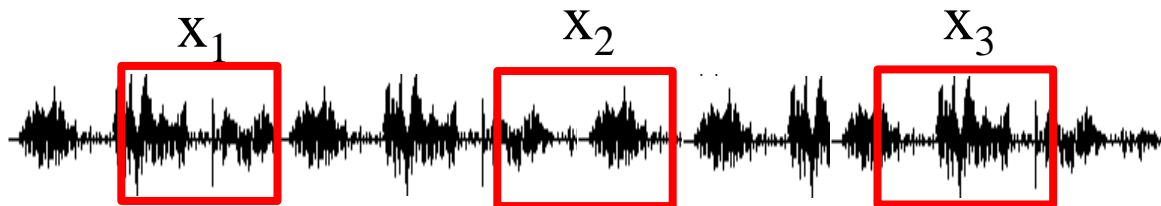


Lattices

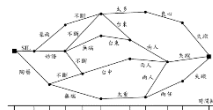


[Chen & Lee,
Interspeech 11]
[Lee & Lee, CSL 14]

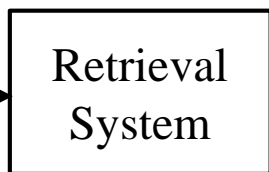
First-pass Retrieval Result



Pseudo Relevance Feedback (PRF)



Lattices

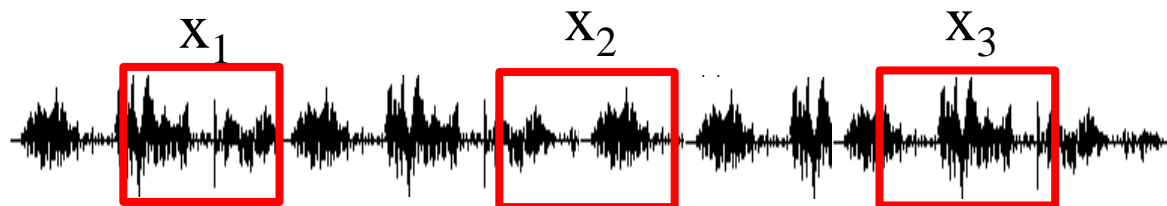


Query Q



[Chen & Lee,
Interspeech 11]
[Lee & Lee, CSL 14]

First-pass Retrieval Result

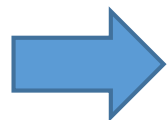


$R(x_1)$

$R(x_2)$

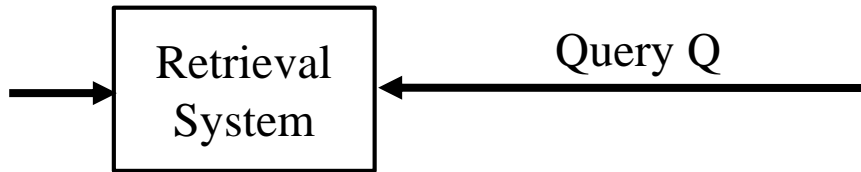
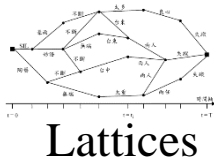
$R(x_3)$

Confidence scores from lattices



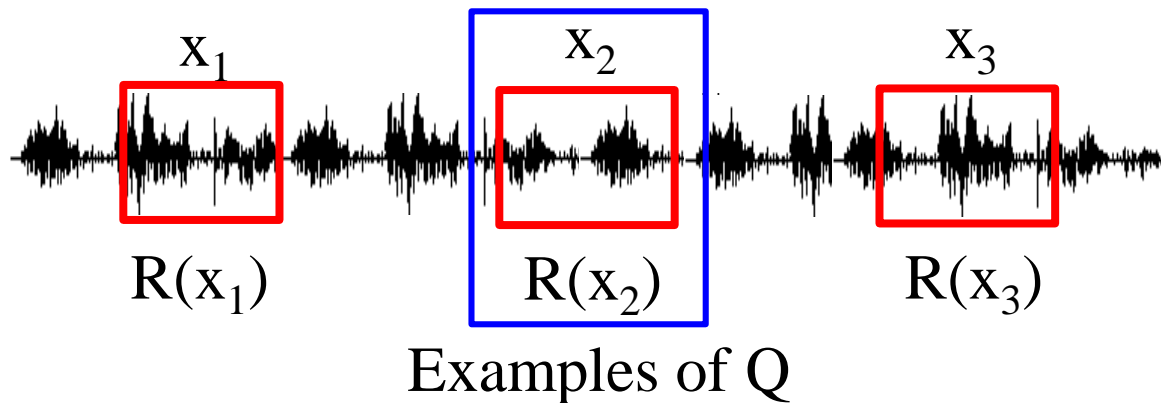
Not shown to the user

Pseudo Relevance Feedback (PRF)



[Chen & Lee,
Interspeech 11]
[Lee & Lee, CSL 14]

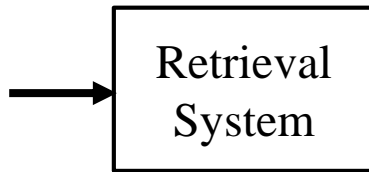
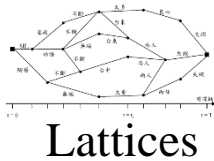
First-pass Retrieval Result



Assume the results with high confidence scores as correct

➡ Considered as examples of Q

Pseudo Relevance Feedback (PRF)

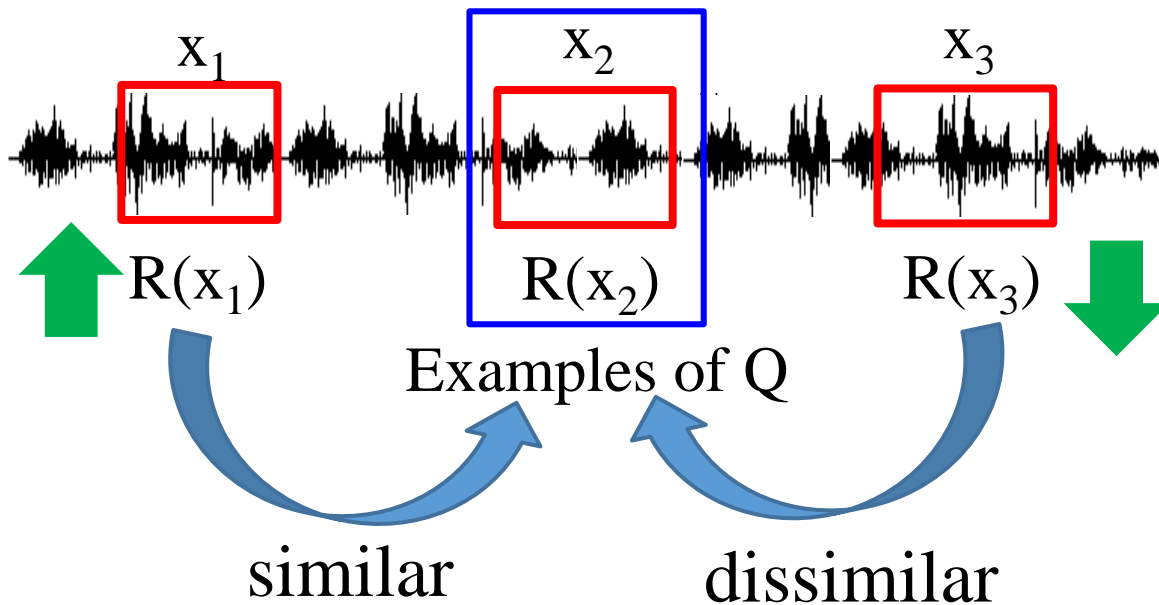


Query Q

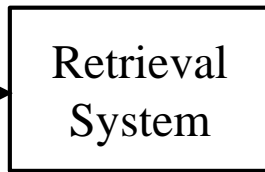
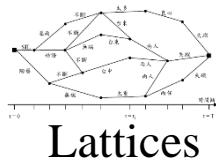


[Chen & Lee,
Interspeech 11]
[Lee & Lee, CSL 14]

First-pass Retrieval Result



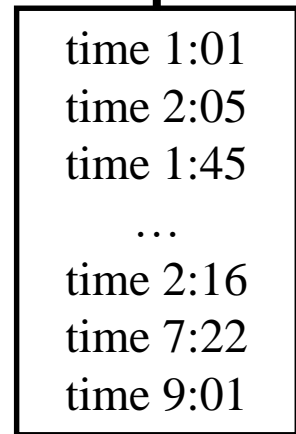
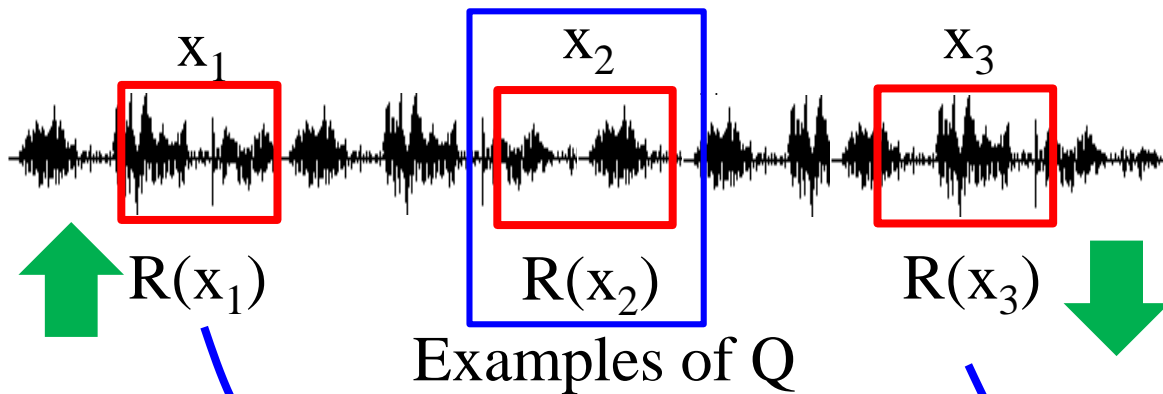
Pseudo Relevance Feedback (PRF)



Query Q



First-pass Retrieval Result



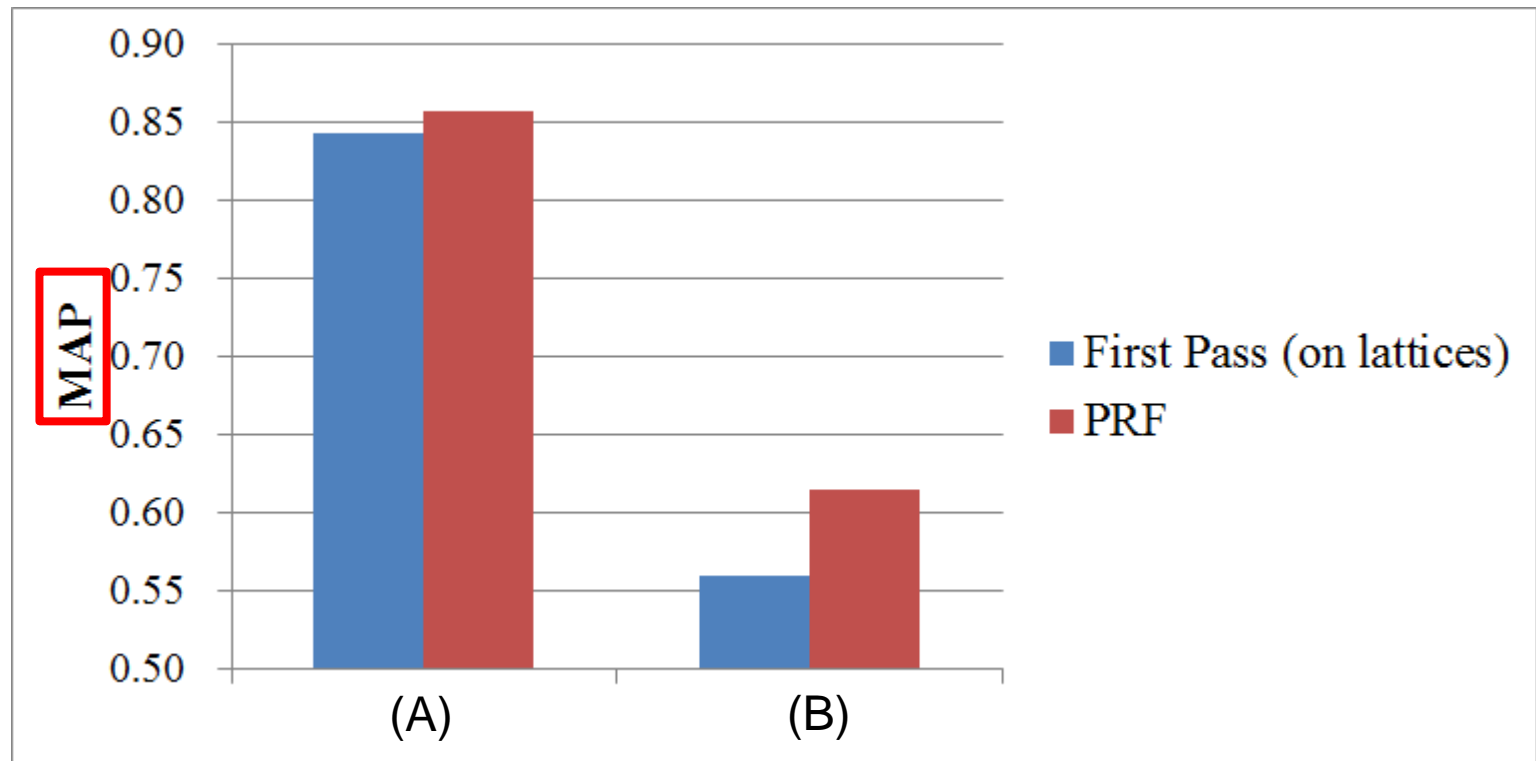
Rank according to new scores

Pseudo Relevance Feedback (PRF)

- Experiments

- Lecture recording [Lee & Lee, CSL 14]

Evaluation Measure: MAP (Mean Average Precision)



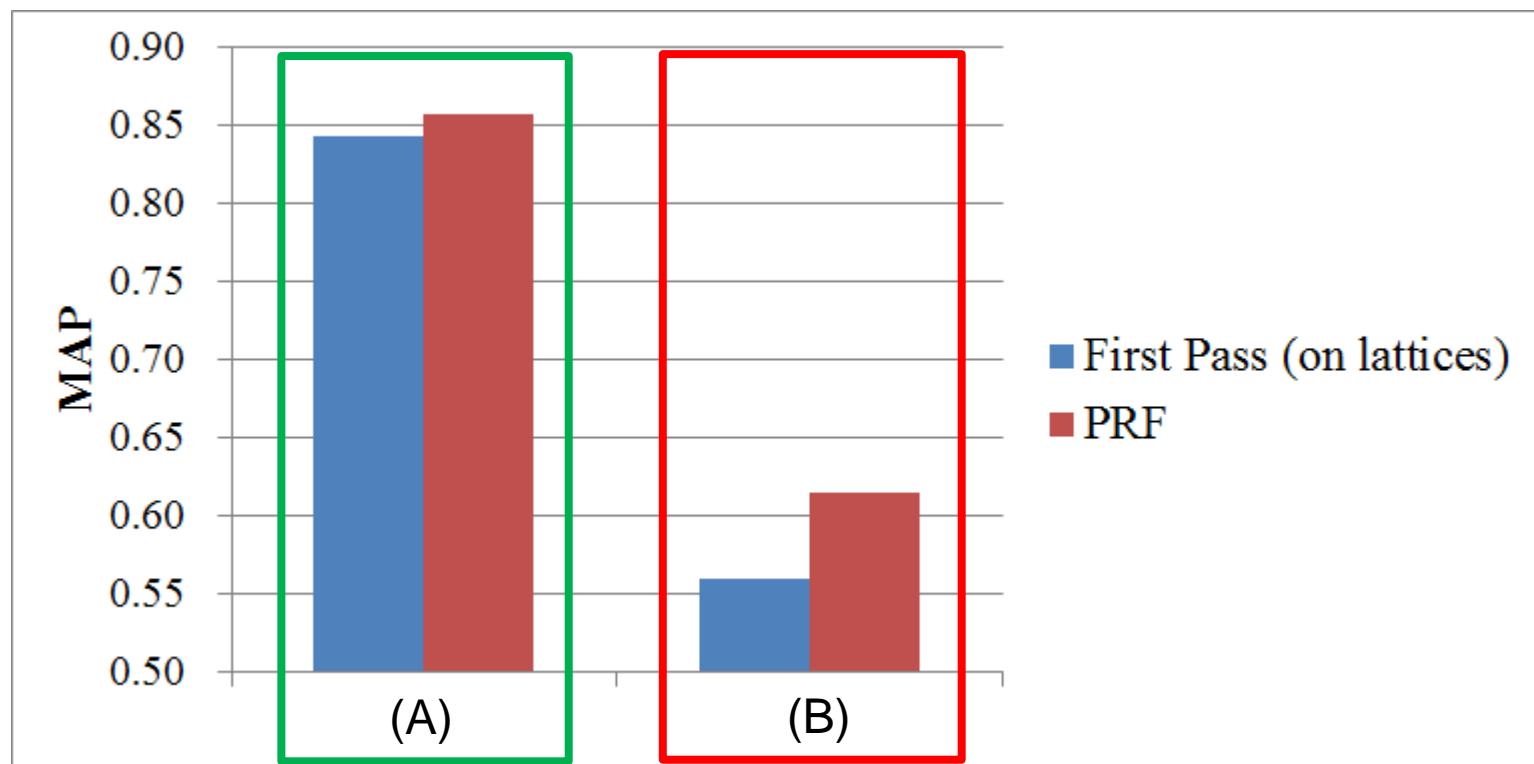
Pseudo Relevance Feedback (PRF)

- Experiments

(A) and (B) use different speech recognition systems

(A): speaker dependent (84% recognition accuracy)

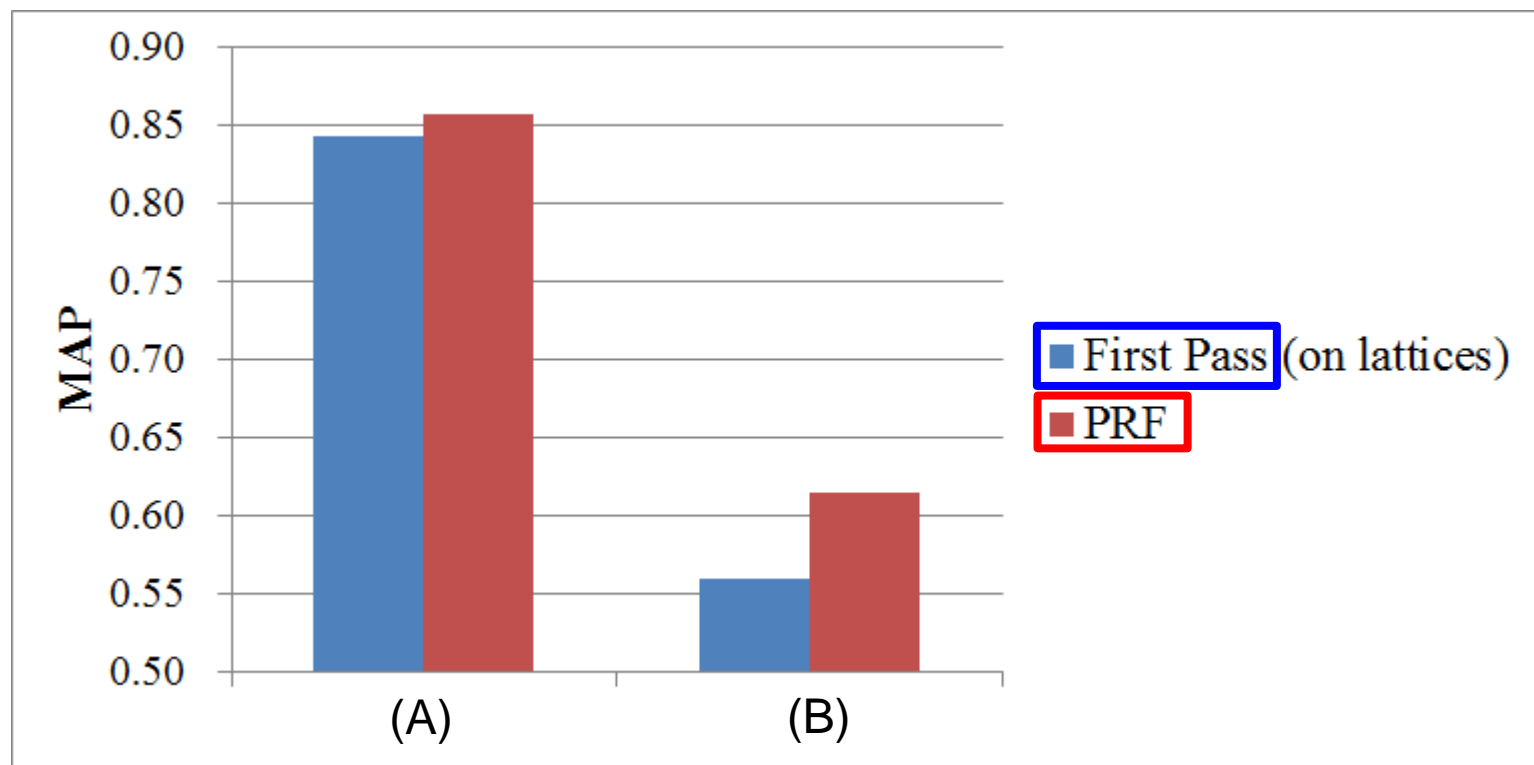
(B): speaker independent (50% recognition accuracy)



Pseudo Relevance Feedback (PRF)

- Experiments

- PRF (red bars) improved the first-pass retrieval results with lattices (blue bars)



New Direction 2-3:
Incorporating
Those Information Lost in ASR
Graph-based Approach



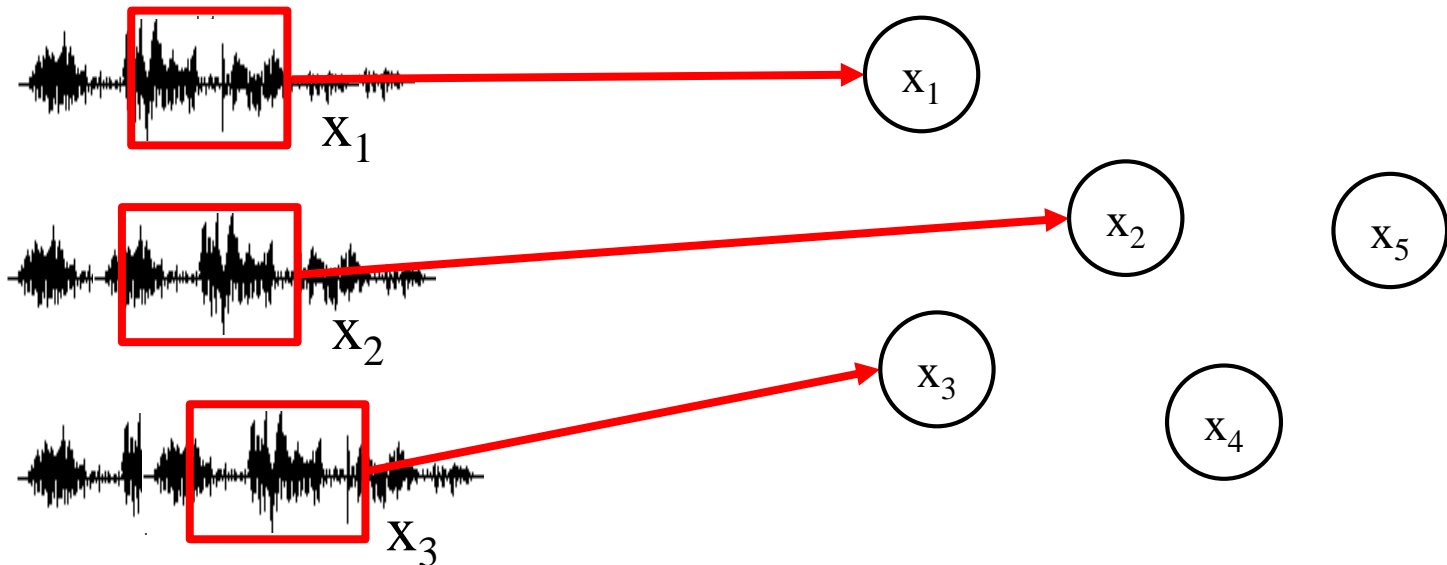
Graph-based Approach

- PRF
 - ▣ Make some assumption to find the examples
 - ▣ Each result considers the similarity to the audio examples
- Graph-based approach [Chen & Lee, ICASSP 11][Lee & Lee, APSIPA 11][Lee & Lee, CSL 14]
 - ▣ Not assume some results are correct
 - ▣ Consider the similarity between all results

Graph Construction

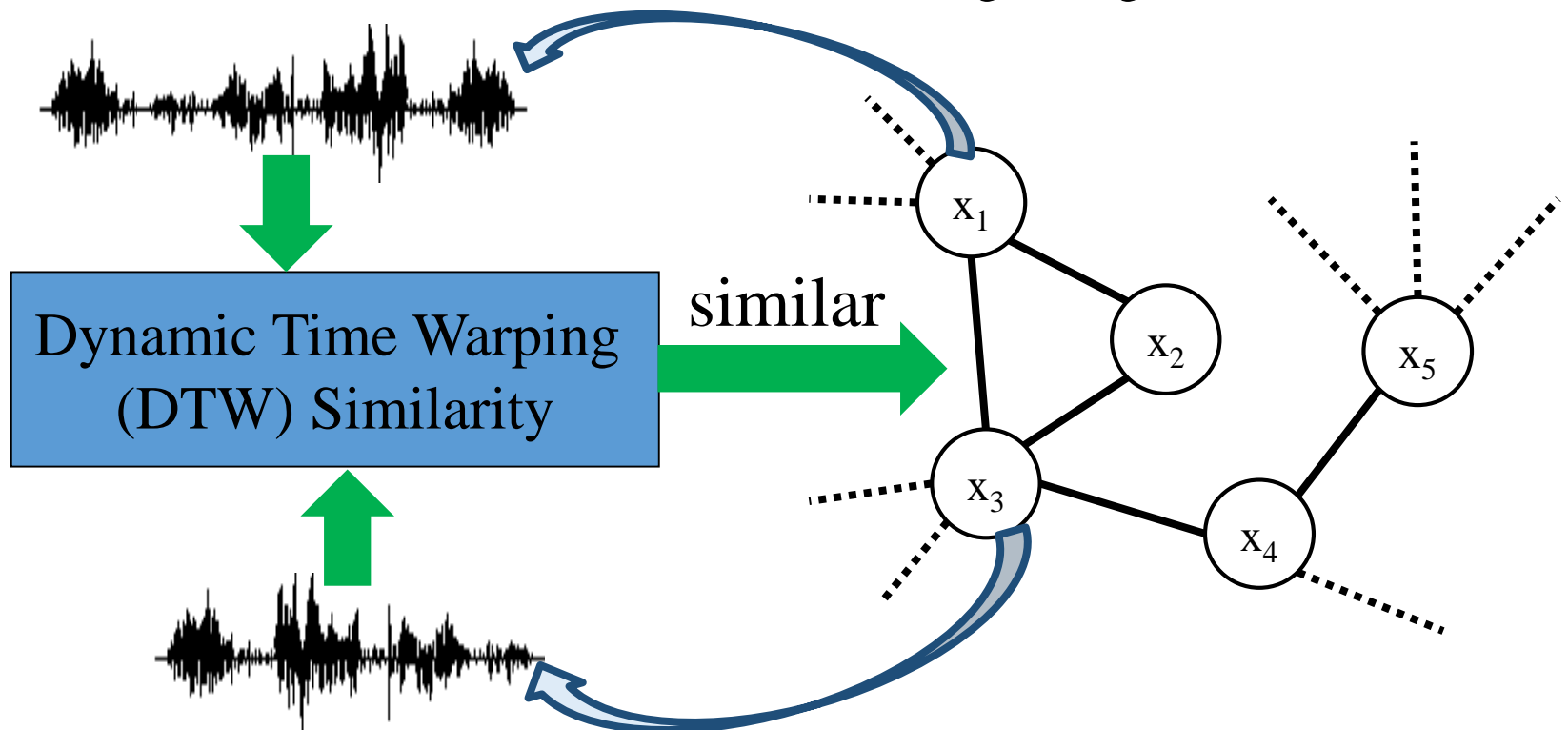
- The first-pass results is considered as a graph.
 - ▣ Each retrieval result is a node

*First-pass Retrieval
Result from lattices*



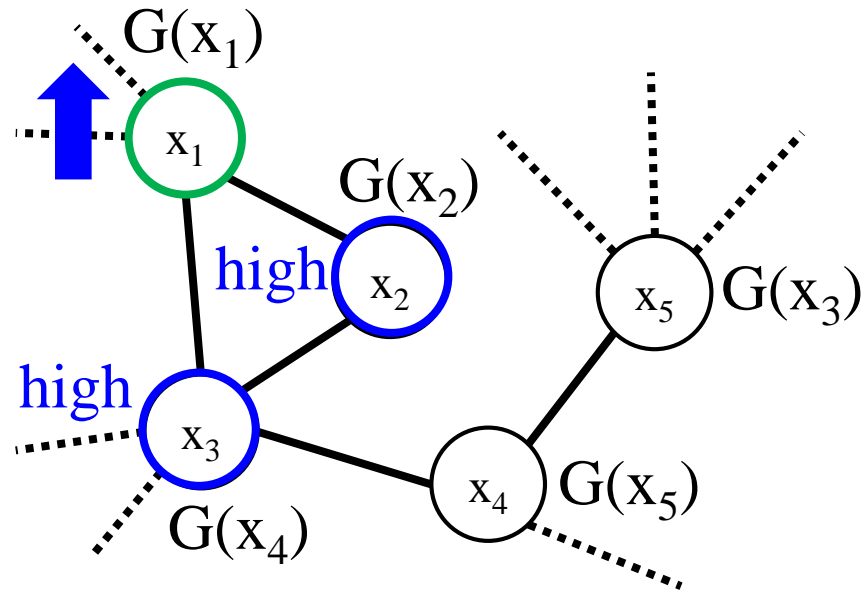
Graph Construction

- The first-pass results is considered as a graph.
 - ▣ Nodes are connected if their retrieval results are similar.
 - DTW similarities are considered as edge weights



Changing Confidence Scores by Graph

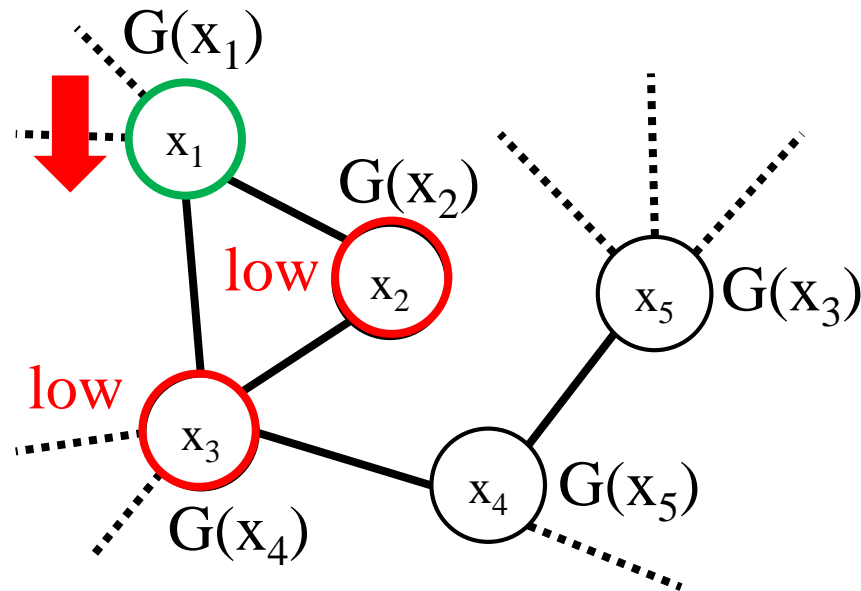
- New scores $G(x_i)$ for each node based on the graph structure **ranked according to new scores**



"You are known by the company you keep"

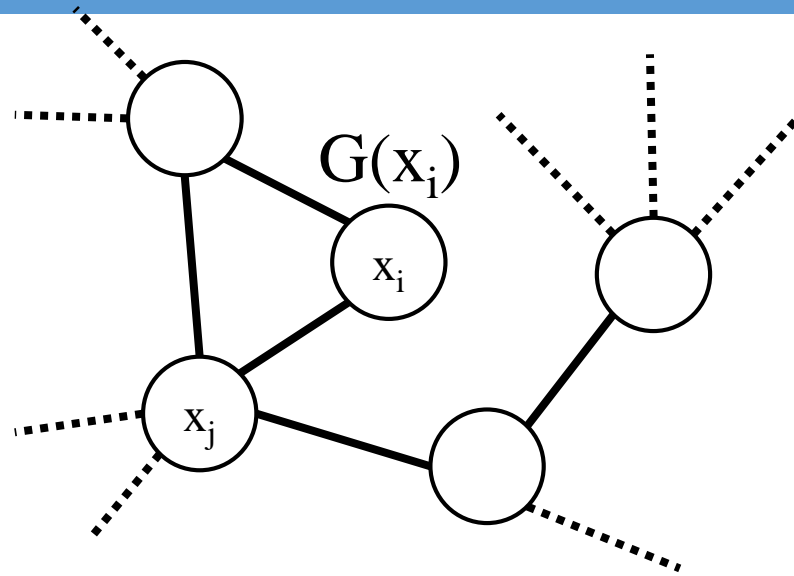
Changing Confidence Scores by Graph

- New scores $G(x_i)$ for each node based on the graph structure **ranked according to new scores**



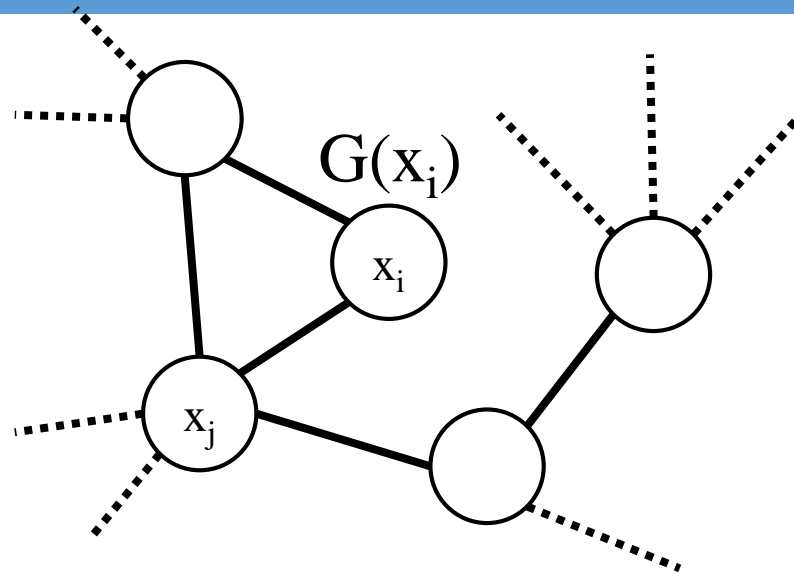
“You are known by the company you keep”

Graph-based Re-ranking - Formulation



$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

Graph-based Re-ranking - Formulation

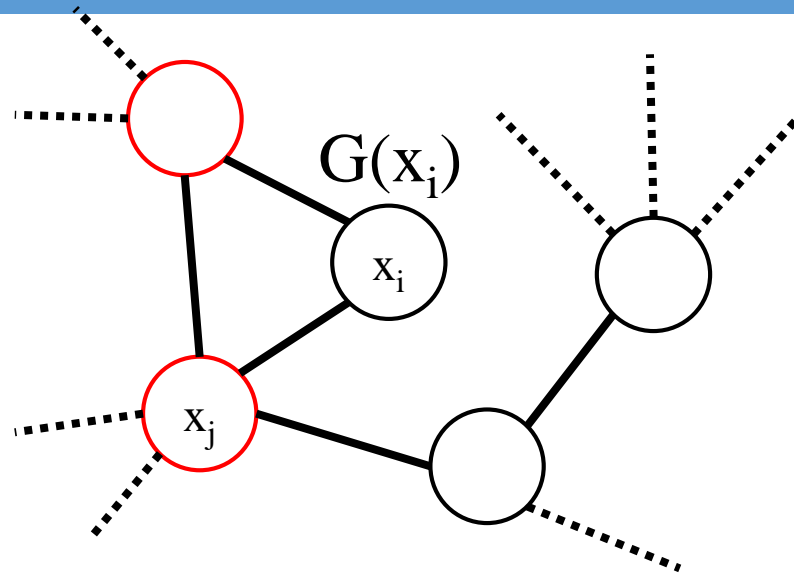


$$G(x_i) = (1 - \alpha) R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

original score
(from lattices)

considering graph structure

Graph-based Re-ranking - Formulation

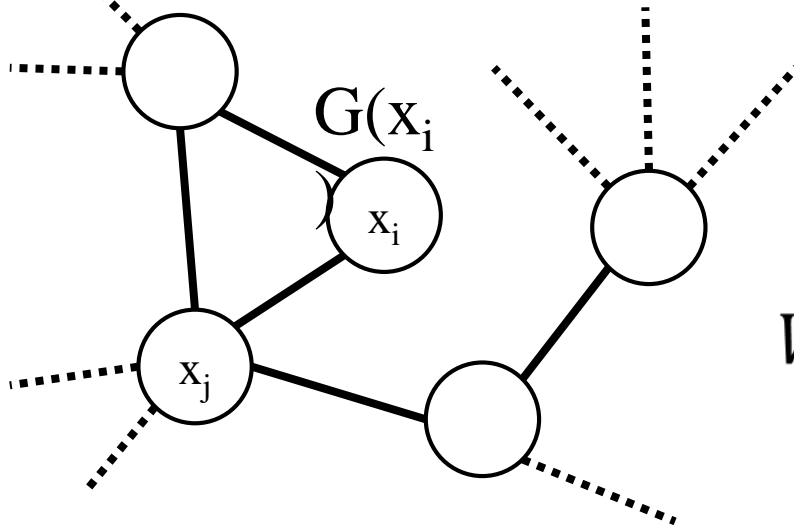


$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

$N(x_i)$: neighbors of x_i (nodes connected to x_i)

x_j : neighbors of x_i (nodes connected to x_i)

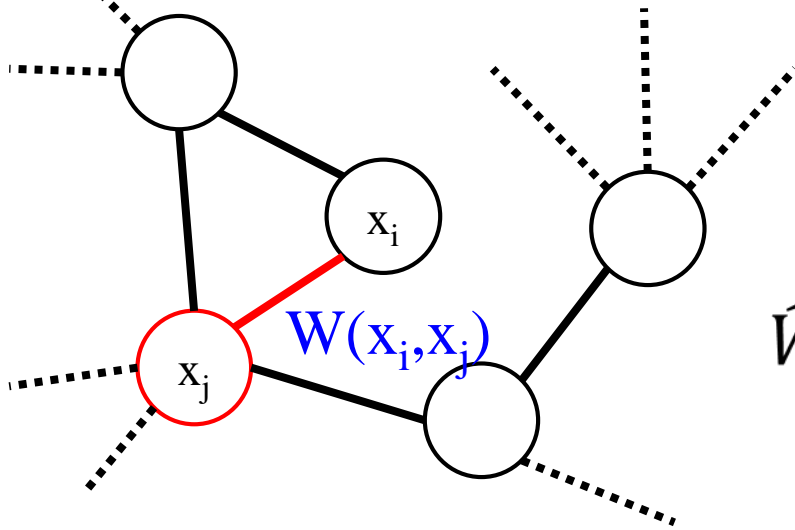
Graph-based Re-ranking - Formulation



$$\hat{W}(x_j, x_i) = \frac{W(x_i, x_j)}{\sum_{x_k \in N(x_j)} W(x_k, x_j)}$$

$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

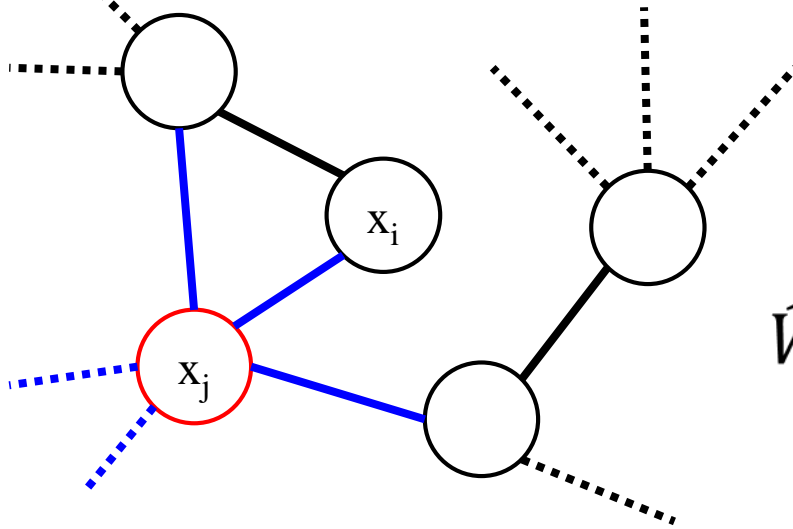
Graph-based Re-ranking - Formulation



$$\hat{W}(x_j, x_i) = \frac{W(x_i, x_j)}{\sum_{x_k \in N(x_j)} W(x_k, x_j)}$$

$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

Graph-based Re-ranking - Formulation

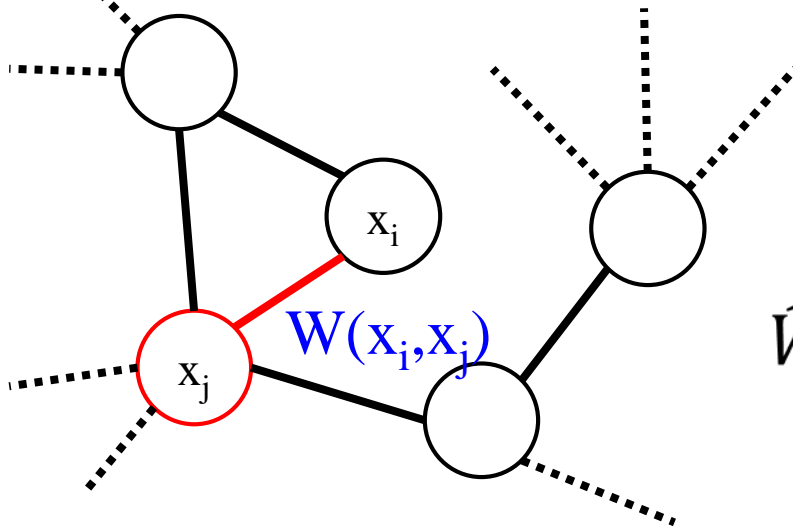


Normalized by the weights of all the edges connected to x_j

$$\hat{W}(x_j, x_i) = \frac{W(x_i, x_j)}{\sum_{x_k \in N(x_j)} W(x_k, x_j)}$$

$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

Graph-based Re-ranking - Formulation



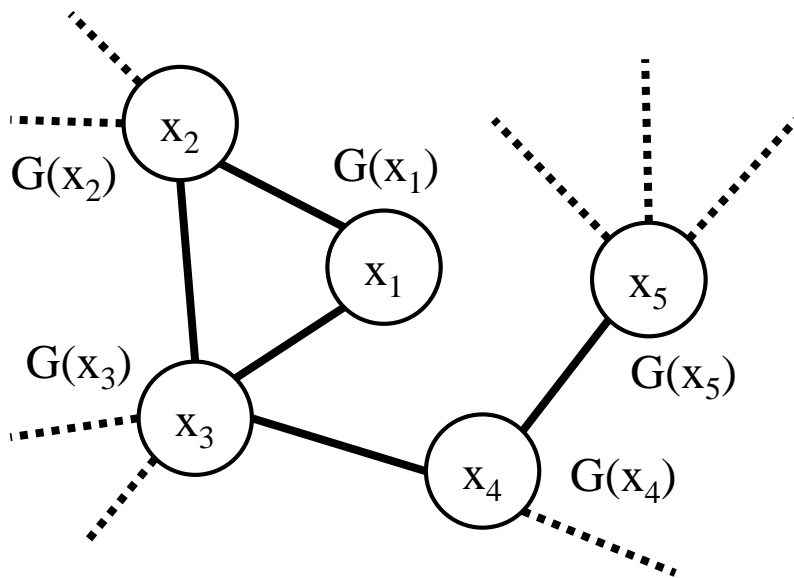
$$\hat{W}(x_j, x_i) = \frac{W(x_i, x_j)}{\sum_{x_k \in N(x_j)} W(x_k, x_j)}$$

$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

The score of x_i would be more close to the nodes x_j with larger edge weights.

Graph-based Re-ranking - Formulation

- Assign score $G(x)$ for each hit region based on the graph structure

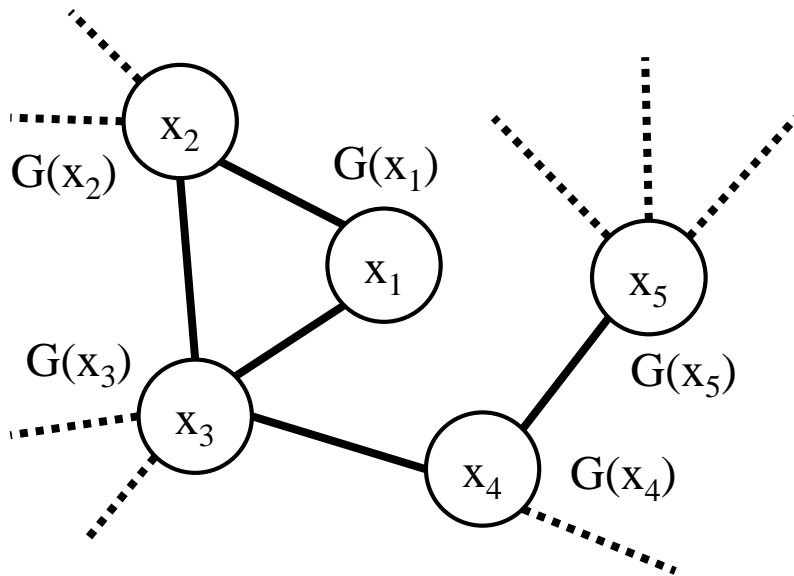


- $G(x_1)$ depends on $G(x_2)$ and $G(x_3)$

$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

Graph-based Re-ranking - Formulation

- Assign score $G(x)$ for each hit region based on the graph structure

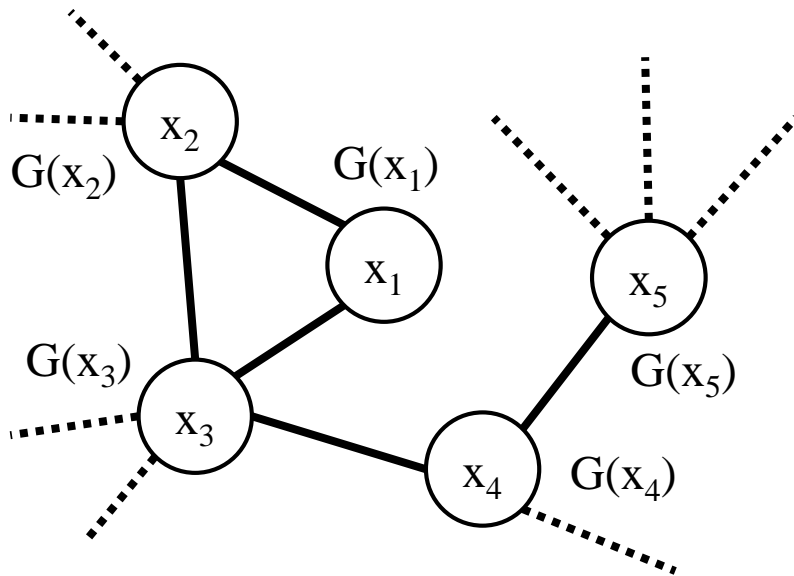


- $G(x_1)$ depends on $G(x_2)$ and $G(x_3)$
- $G(x_2)$ depends on $G(x_1)$ and $G(x_3)$

$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

Graph-based Re-ranking - Formulation

- Assign score $G(x)$ for each hit region based on the graph structure

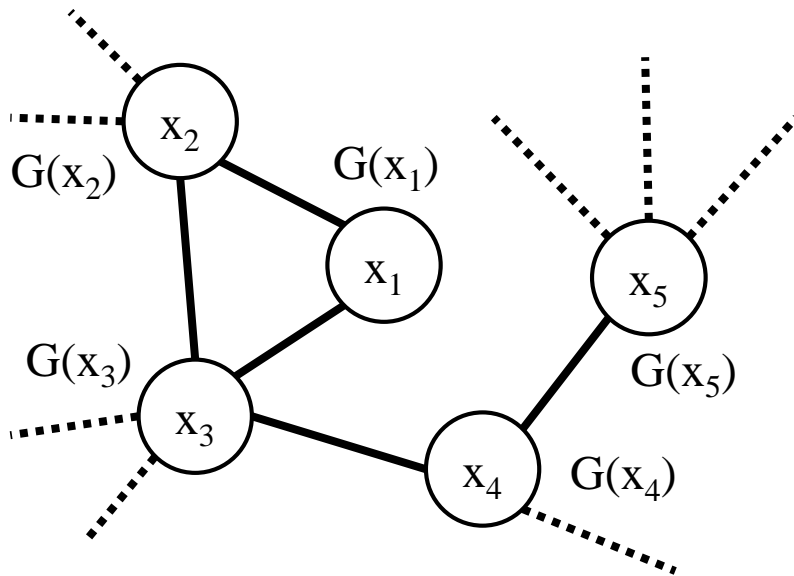


- $G(x_1)$ depends on $G(x_2)$ and $G(x_3)$
- $G(x_2)$ depends on $G(x_1)$ and $G(x_3)$
-

$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

Graph-based Re-ranking - Formulation

- Assign score $G(x)$ for each hit region based on the graph structure



- How to find $G(x_1), G(x_2), G(x_3) \dots$ satisfying the following equation?

- This is random walk.

$G(x_i)$ is uniquely and efficiently obtainable

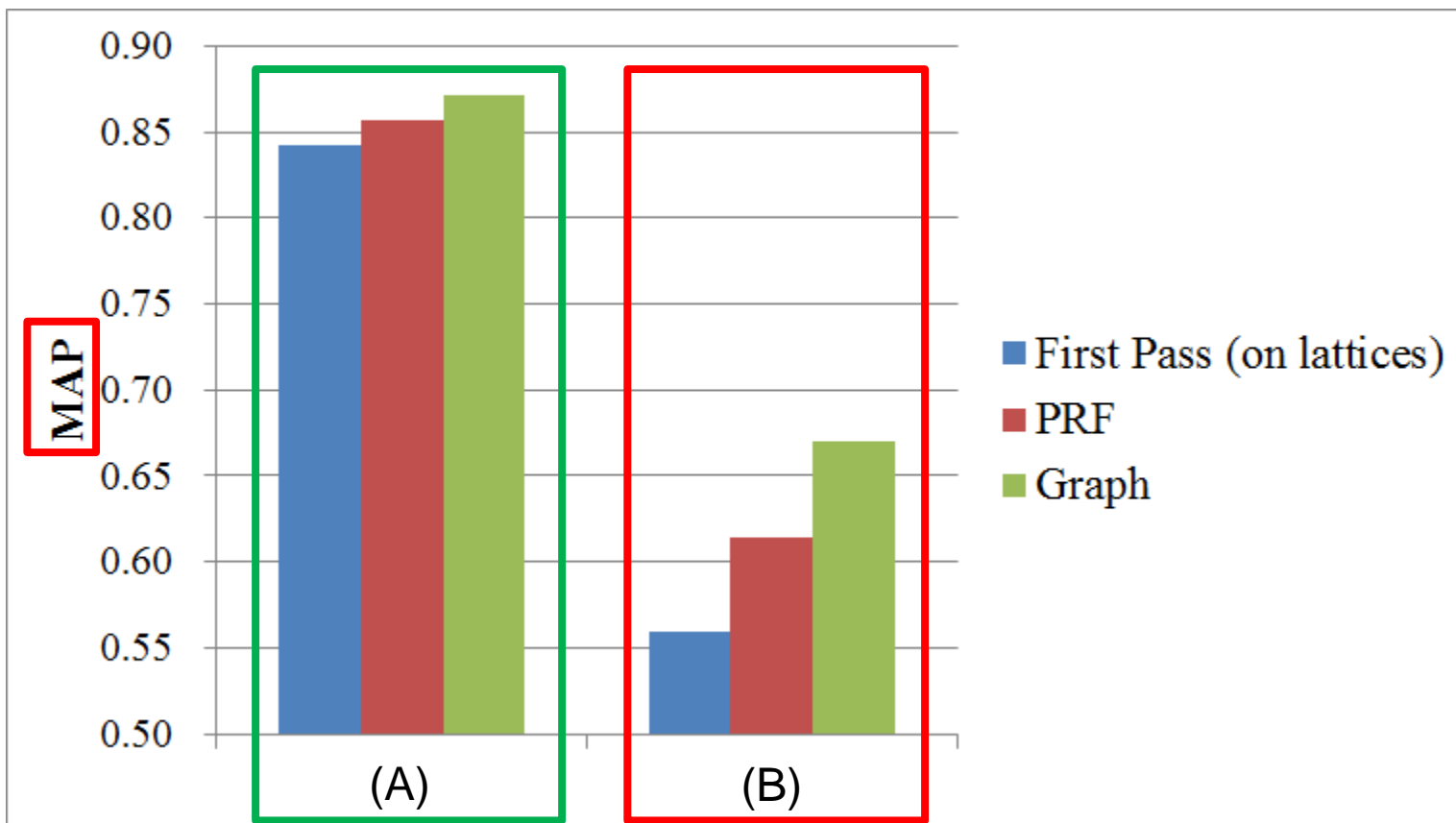
$$G(x_i) = (1 - \alpha)R(x_i) + \alpha \sum_{x_j \in N(x_i)} G(x_j) \hat{W}(x_j, x_i)$$

Graph-based Approach - Experiments

- Lecture recording [Lee & Lee, CSL 14]

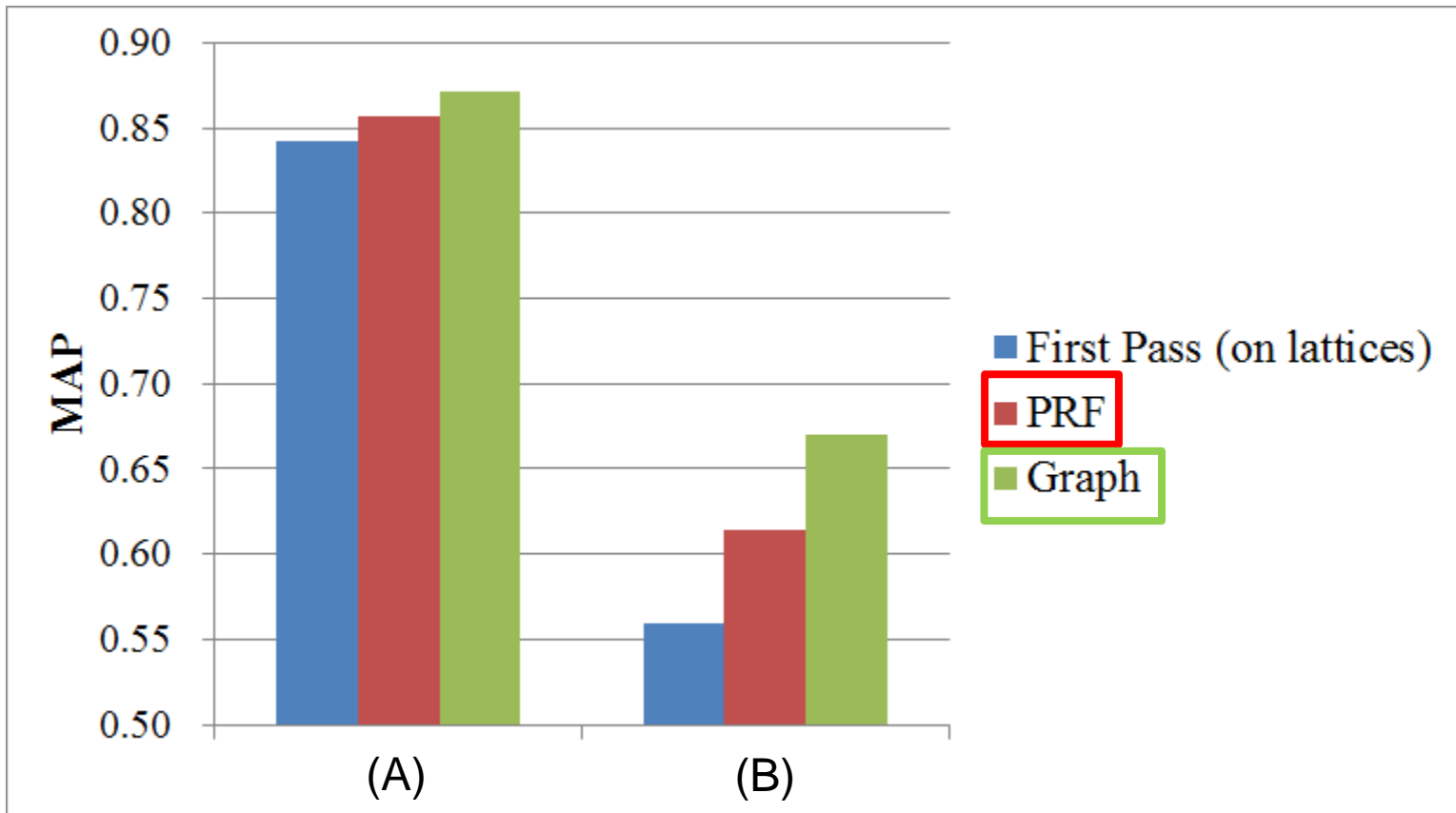
(A): speaker dependent (high recognition accuracy)

(B): speaker independent (low recognition accuracy)



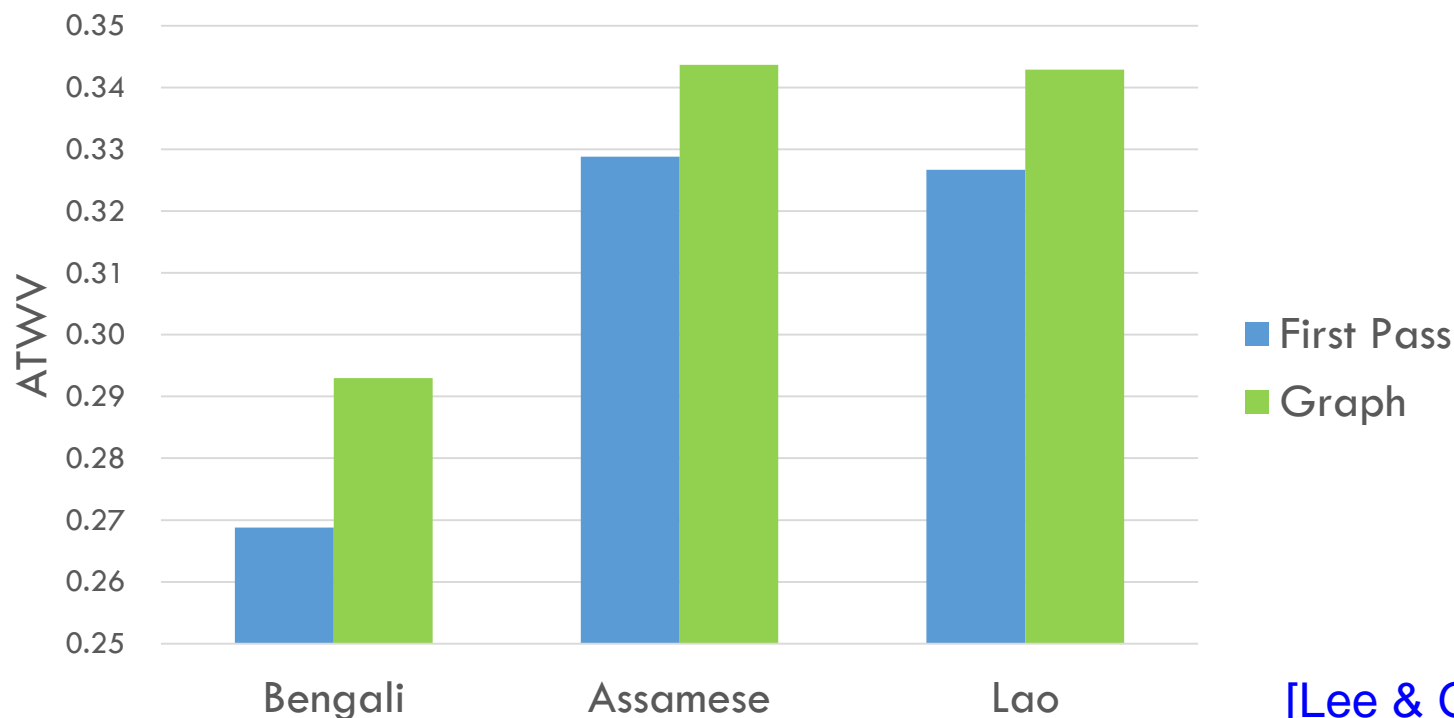
Graph-based Approach - Experiments

- Graph-based re-ranking (green bars) outperformed PRF (red bars)



Graph-based Approach – Experiments

- Graph-based approach on limited language data from the IARPA Babel program



[Lee & Glass,
Interspeech 14]

Graph-based Approach – Experiments

- 13% relative improvement on OOV queries on lecture recording (several speakers) [Jansen, ICASSP 13][Norouzian, ICASSP 13]
- 14% relative improvement on AMI Meeting Corpus [Norouzian, Interspeech 13]
 - ▣ Graph Spectral Clustering
- Optimizing evaluation measure and considering the graph structure at the same time [Audhkhasi, ICASSP 2014]
- 11% relative improvement with subword-based system on OpenKWS15 (Swahili) [Van Tung Pham, ICASSP, 2016]

New Direction 3:
No Speech Recognition!



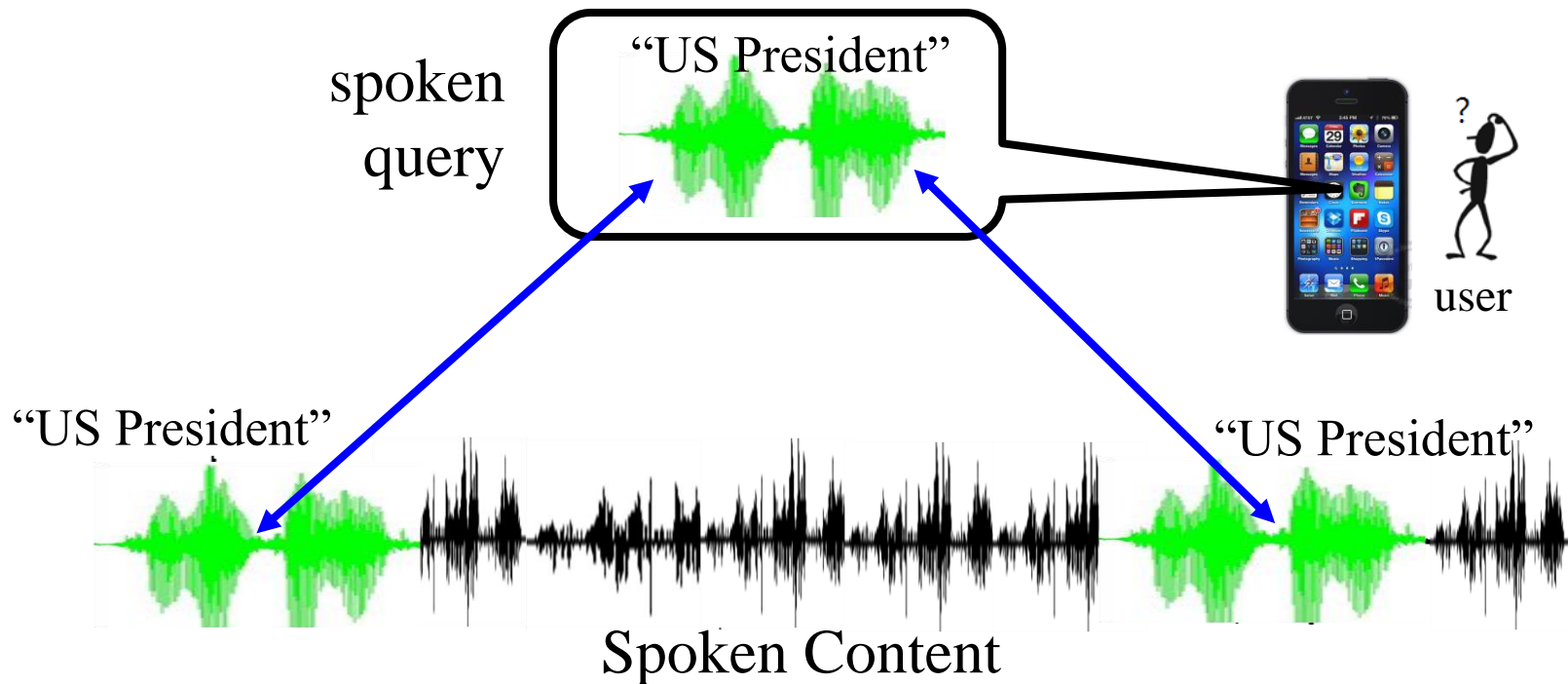
Why Spoken Content Retrieval without Speech Recognition?

- Bypassing ASR to avoid information loss and all problems with ASR (errors, OOV words, background noise, etc.)
- Just to identify the query, no need to find out which words the query includes
- Audio files on the Internet in hundreds of different languages
 - ▣ Too limited annotated data for training reliable speech recognition systems for most languages
 - ▣ Written form even doesn't exist for some languages
- Many audio files are code-switched across several different languages

Spoken Content Retrieval without Speech Recognition

Is it possible?

This task is called “Query-by-Example”.



Compute similarity between spoken queries and audio files on acoustic level, and find the query term

Approach Categories

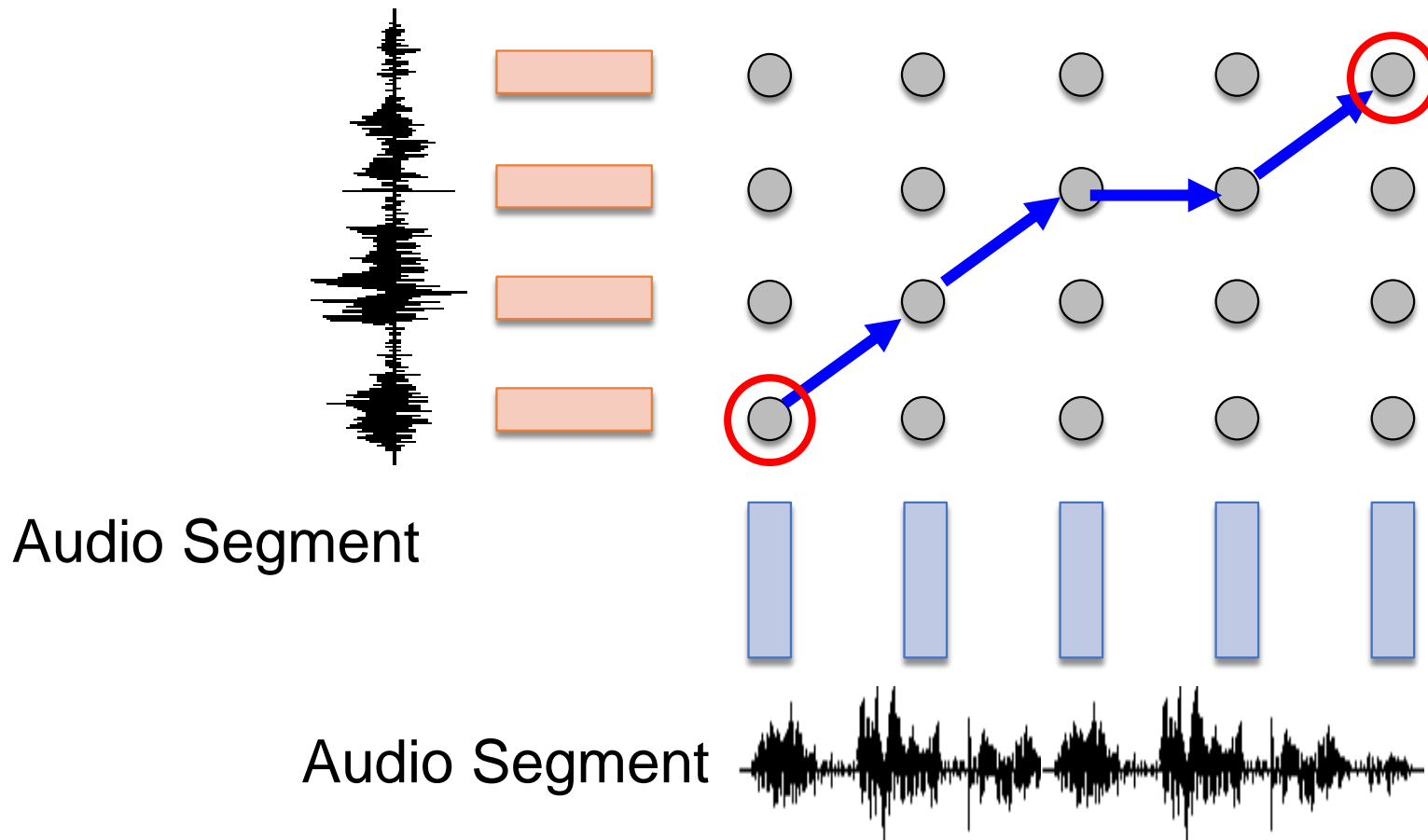
- DTW-based Approaches
 - ▣ Matching sequences with DTW
- Audio Segment Representation
 - ▣ Representing audio segments by fixed length vector representations
- Unsupervised ASR (or model-based approach)
 - ▣ Training word- or subword-like acoustic patterns (or tokens) from target audio archive
 - ▣ Transcribing both the audio archive and the query into word- or subword-like token sequences
 - ▣ Matching based on the tokens, just like text retrieval

New Direction 3-1:
No Speech Recognition!
DTW-based Approaches



DTW-based Approach

□ Conventional DTW

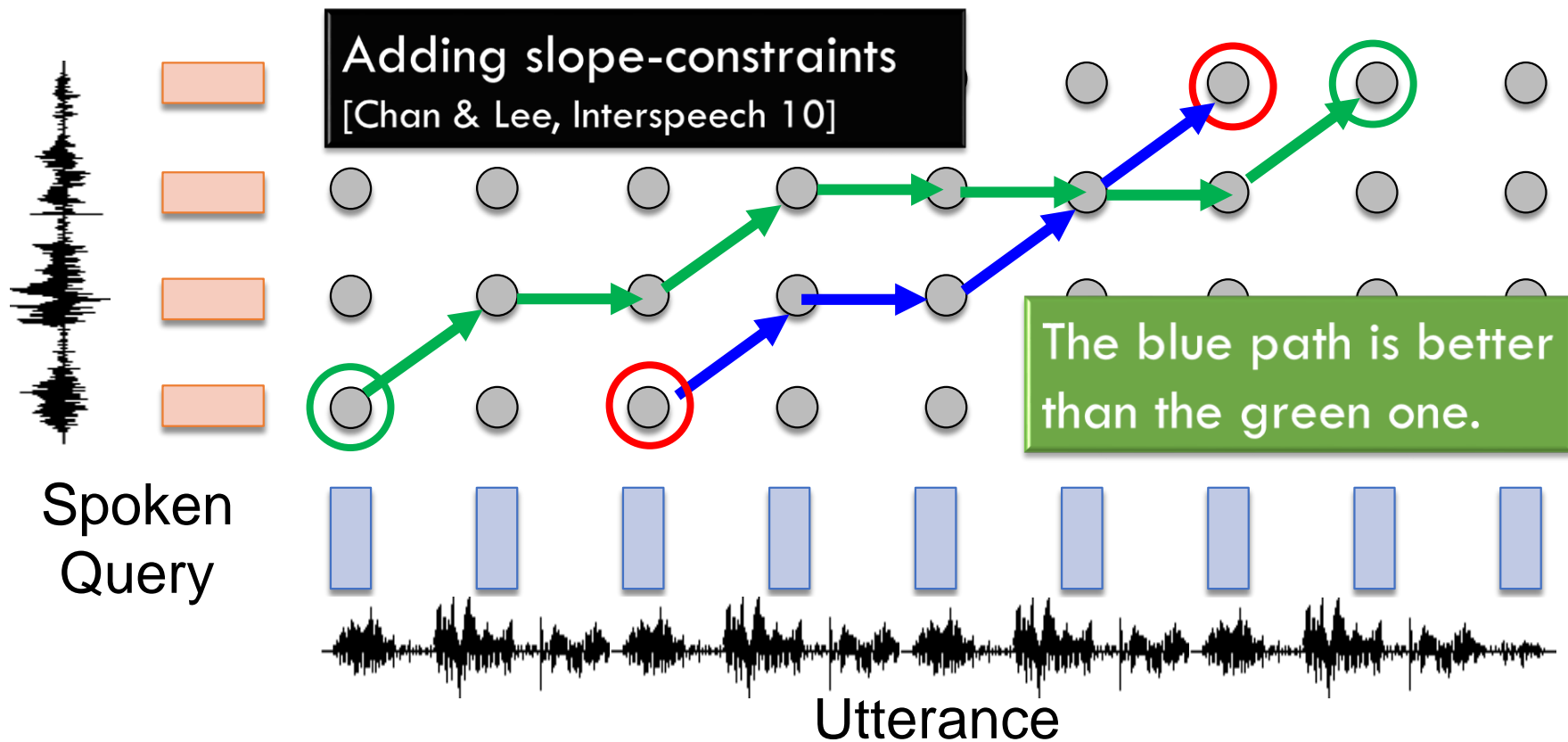


DTW-based Approach

- DTW for query-by-example

- Whether a spoken query is in an utterance

Segmental DTW [Zhang, ICASSP 10], Subsequence DTW [Anguera, ICME 13][Calvo, MediaEval 14]

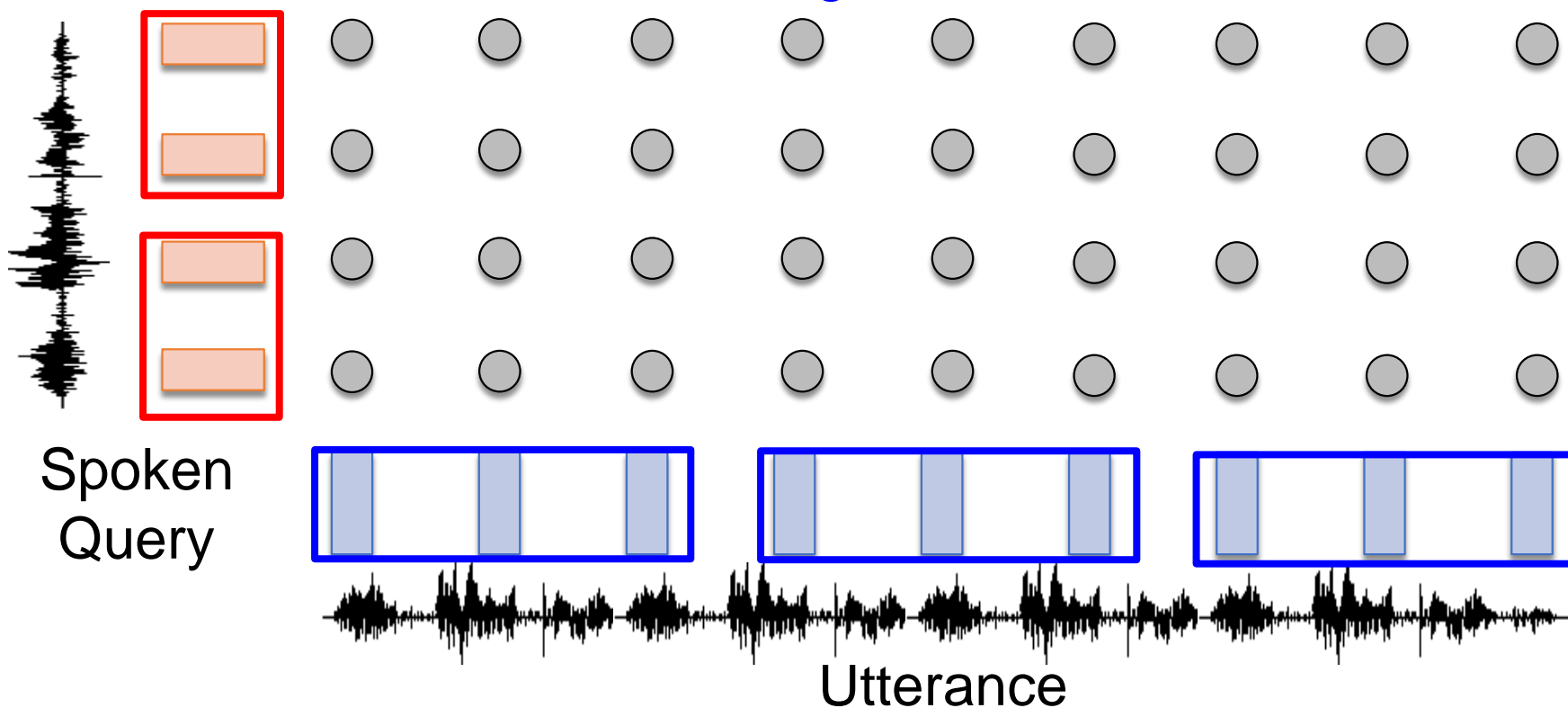


Acoustic Feature Vectors

- Gaussian posteriorgram [Zhang, ICASSP 10][Wang, MediaEval 14]
- Phonetic posteriors [Hazen, ASRU 09]
 - ▣ MLP trained from another corpus (probably in a different language)
- Bottle-neck feature generated from MLP [Kesiraju, MediaEval 14]
- RBM posteriorgram [Zhang, ICASSP 12]
- Performance comparison [Carlin, Interspeech 11]

Speed-up Approaches for DTW

- Segment-based matching [Chan & Lee, Interspeech 10][Chan & Lee, ICASSP 11] Group consecutive acoustically similar feature vectors into a segment



Speed-up Approaches for DTW

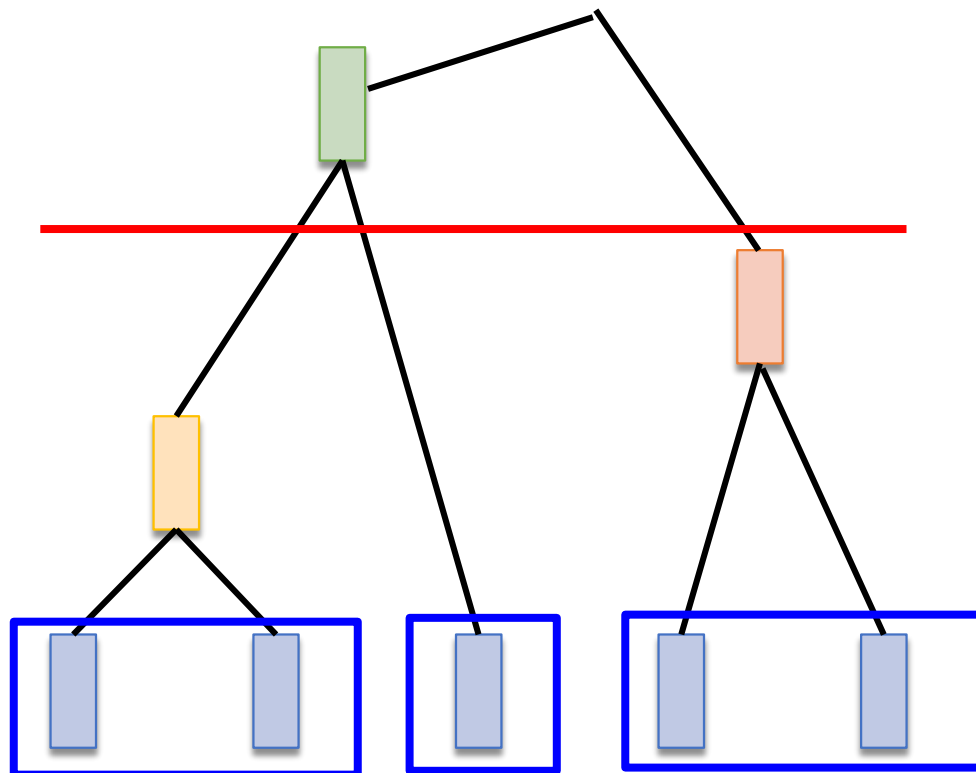
- Segment-based matching

Group consecutive acoustically similar feature vectors into a segment

Hierarchical
Agglomerative
Clustering (HAC)

Step 1: build a tree

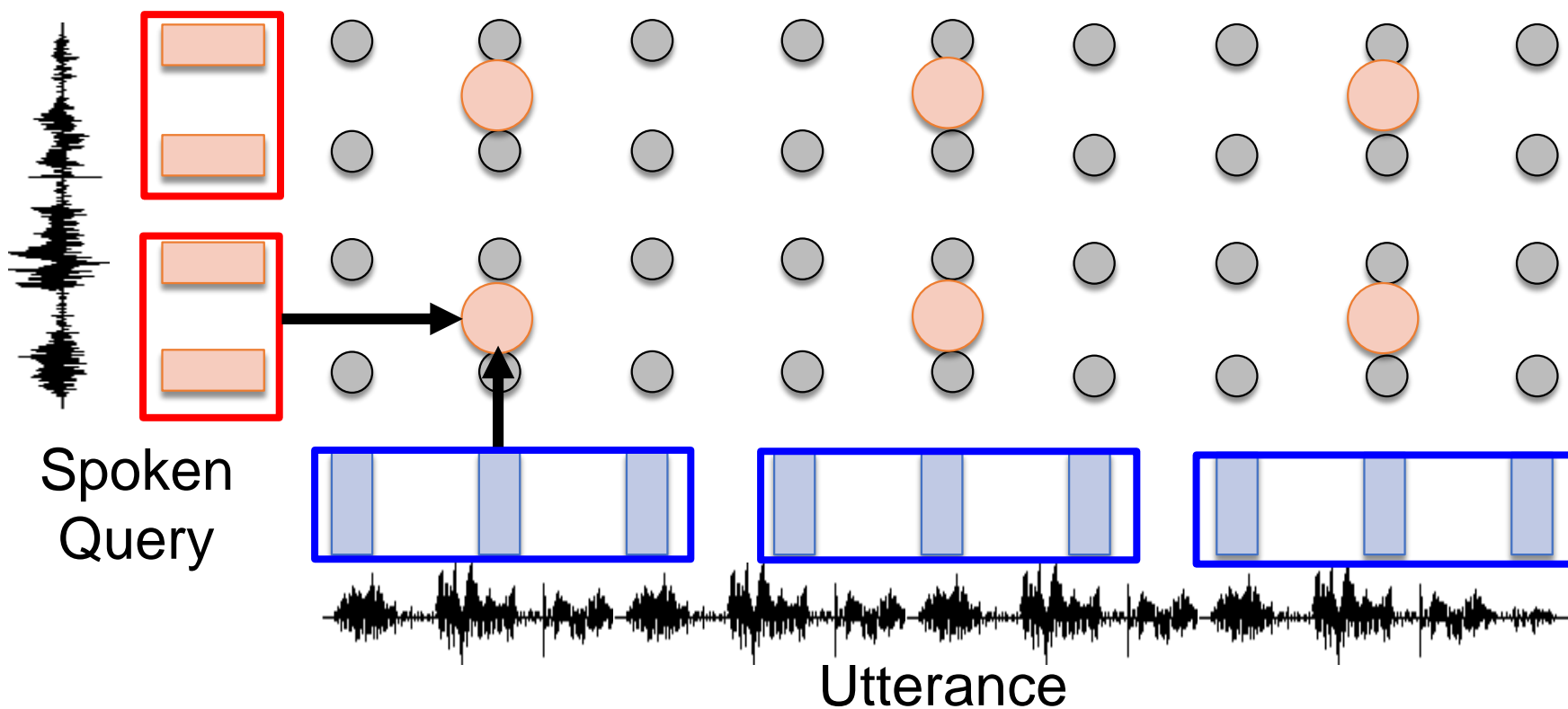
Step 2: pick a
threshold



Speed-up Approaches for DTW

- Segment-based matching [Chan & Lee, Interspeech 10][Chan & Lee, ICASSP 11]

Compute similarities between segments only



Speed-up Approaches for DTW

- Segment-based matching [Chan & Lee, Interspeech 10][Chan & Lee, ICASSP 11]
- Lower bound estimation [Zhang, ICASSP 11][Zhang, Interspeech 11]
- Indexing the frames in the target audio file [Jansen, ASRU 11][Jansen, Interspeech 12]
- Information Retrieval based DTW [Anguera, Interspeech 13]

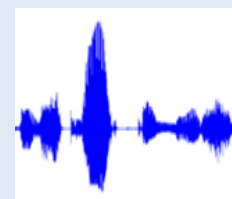
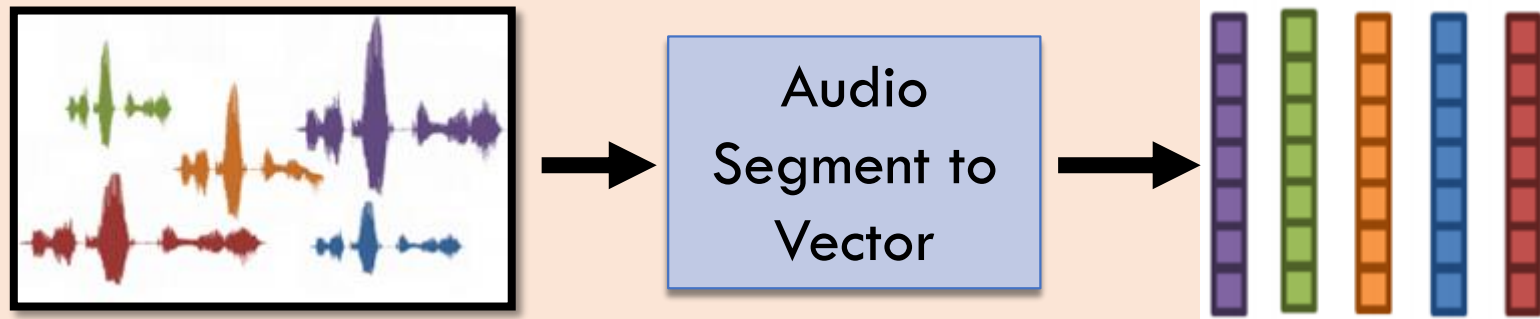
New Direction 3-2:
No Speech Recognition!
Audio Segment Representation

Framework

[Chung & Lee, Interspeech 16][Chen, ICASSP 15]
[Levin, ICASSP 15][Levin, ASRU 13]

Audio archive divided into variable-length audio segments

Off-line



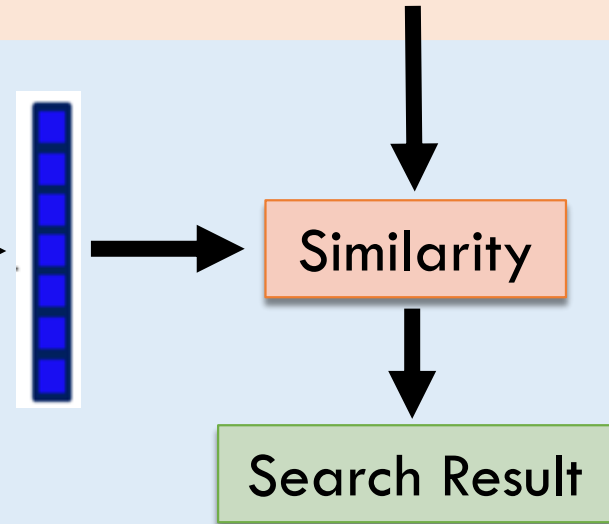
Spoken Query

On-line



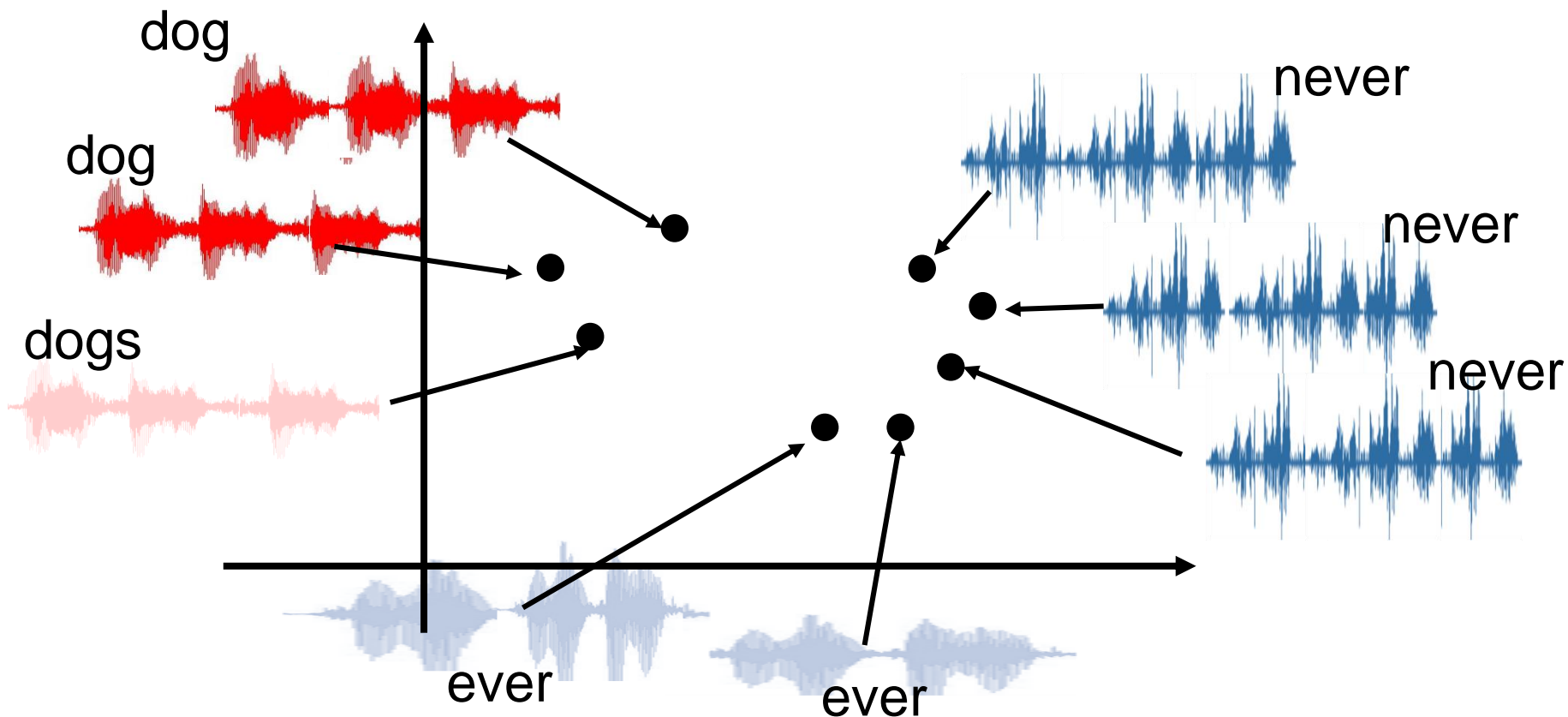
Similarity

Search Result



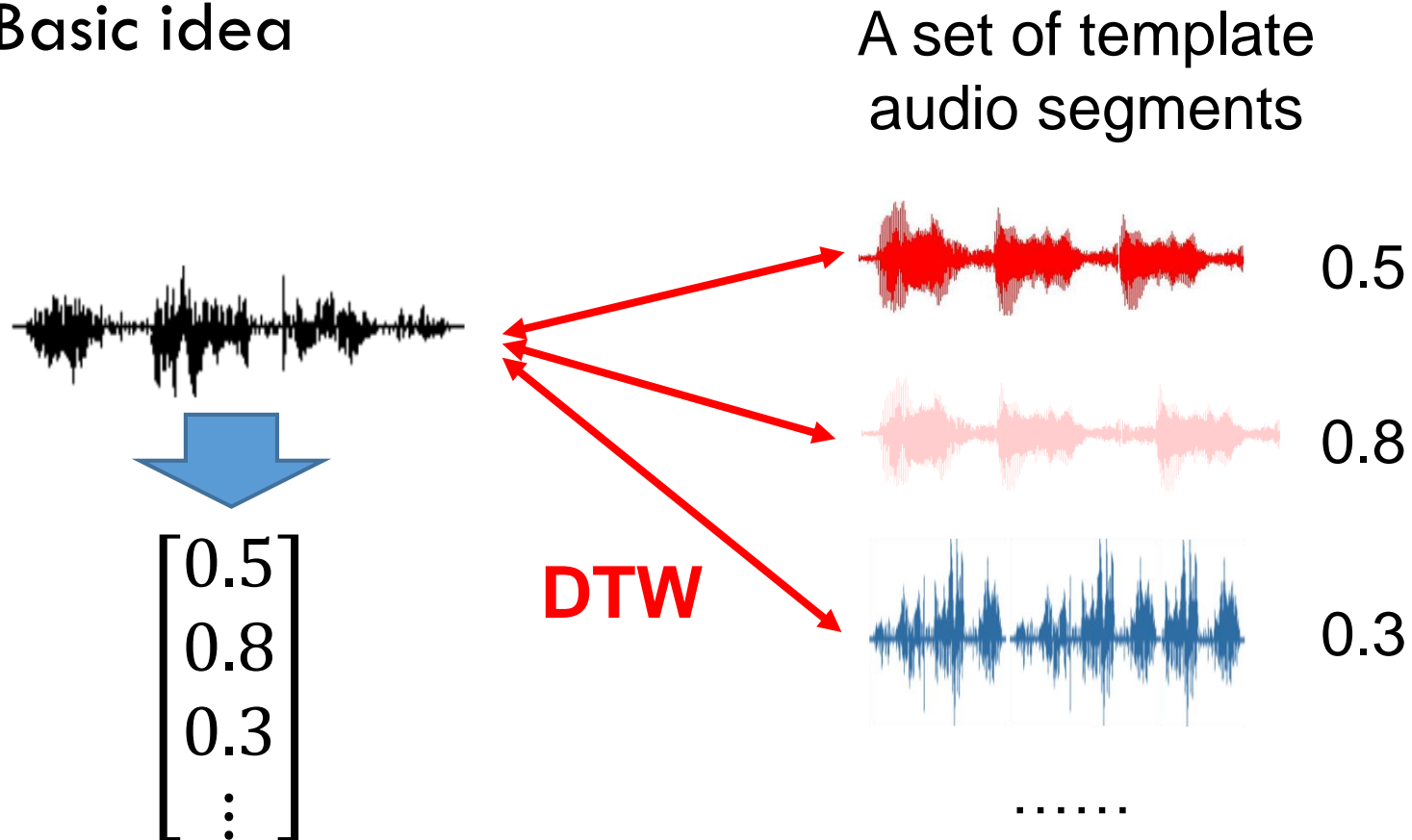
Audio Word to Vector

- The audio segments corresponding to words with similar pronunciations are close to each other.



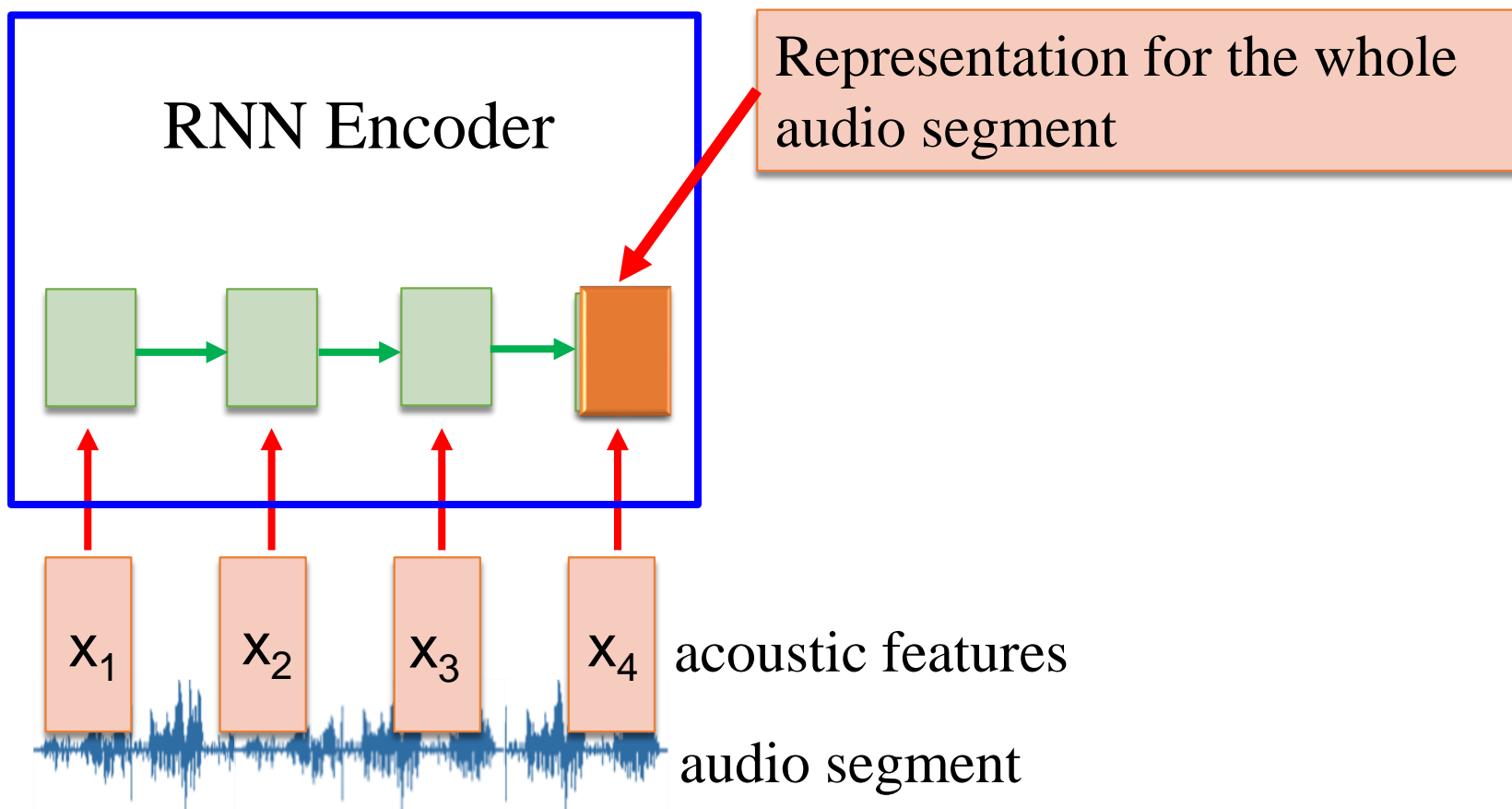
Audio Word to Vector - Segmental Acoustic Indexing

□ Basic idea



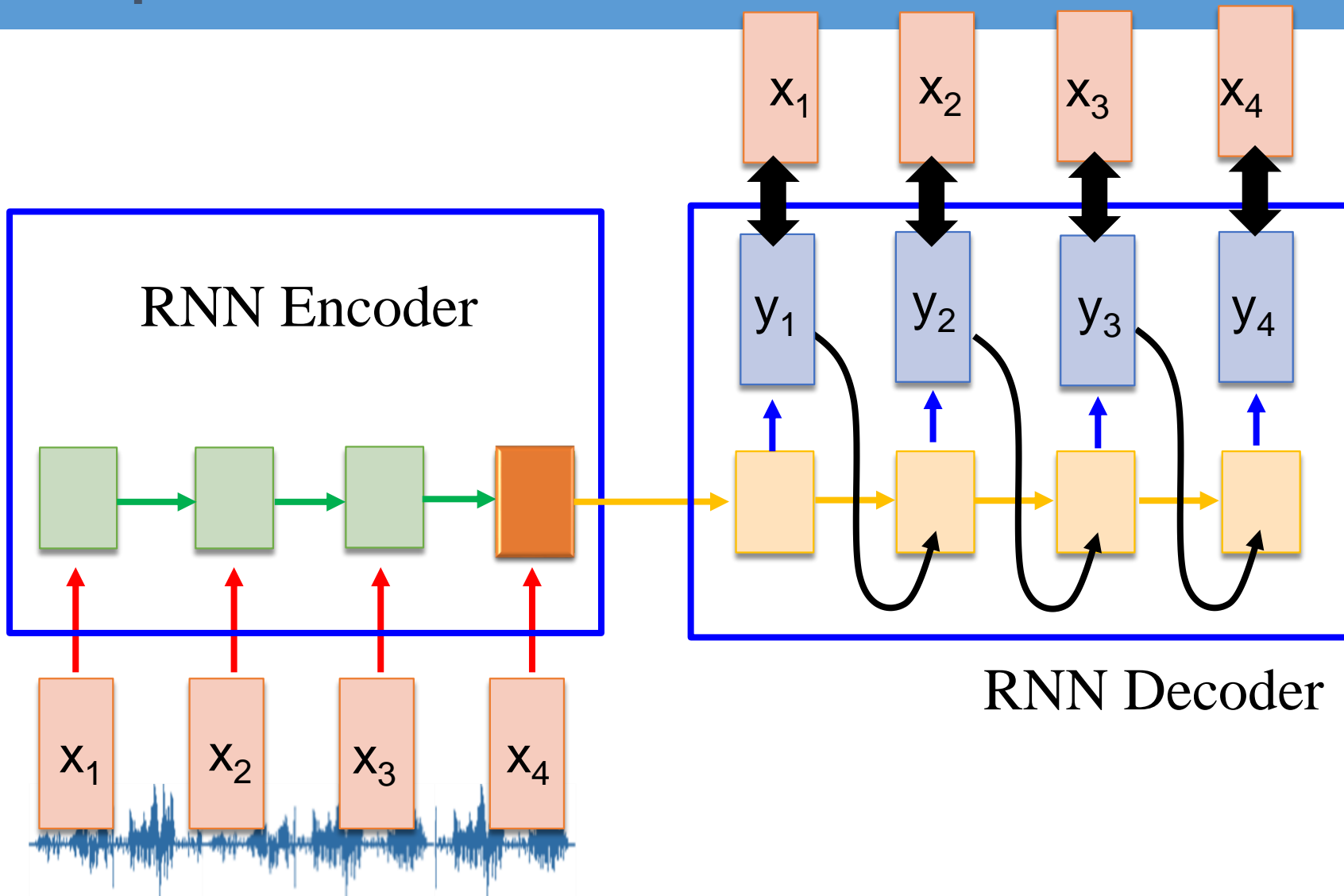
Audio Word to Vector – Sequence Auto-encoder

[Chung & Lee,
Interspeech 16]

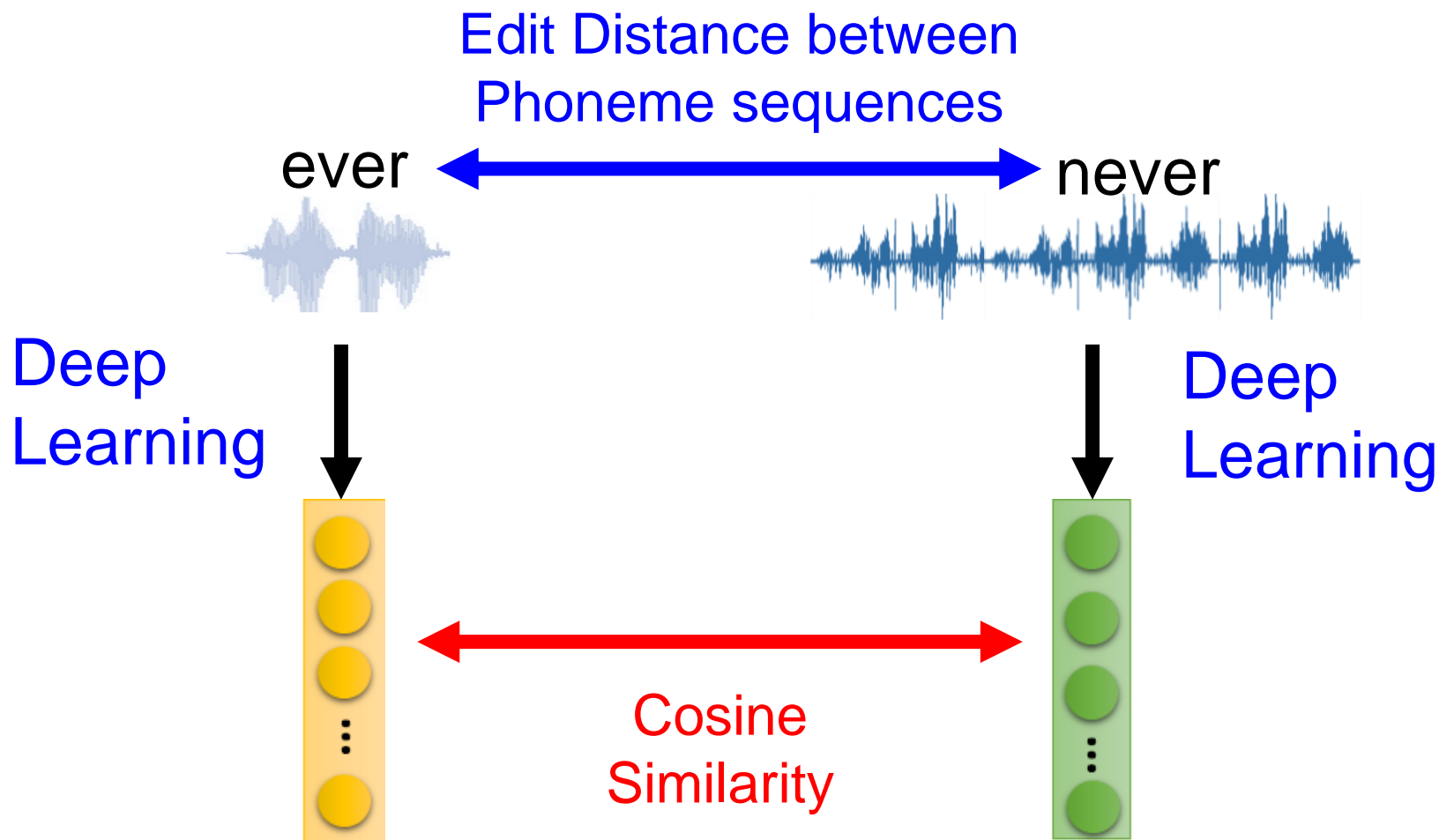


Audio Word to Vector – Sequence Auto-encoder

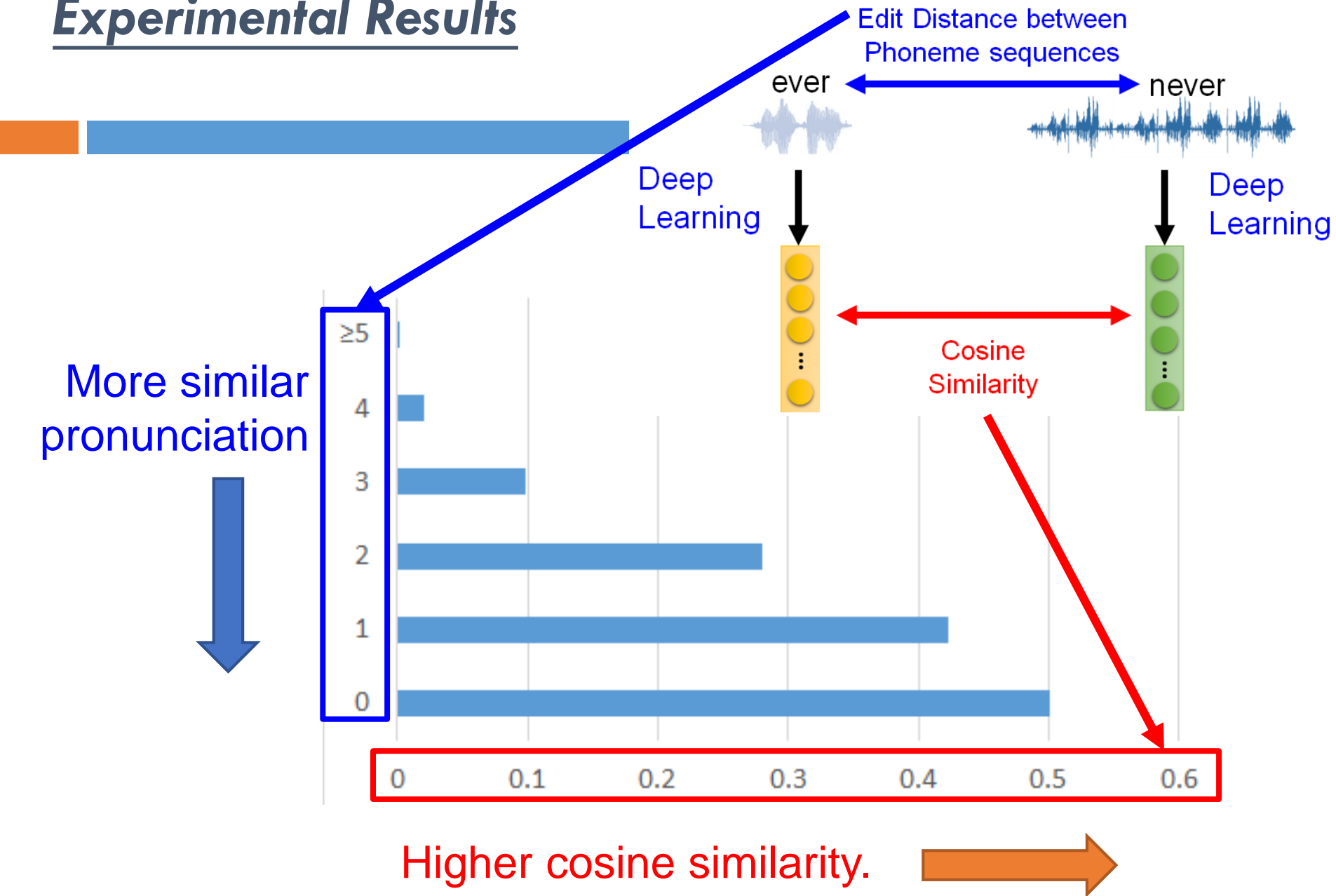
[Chung & Lee,
Interspeech 16]



Sequence Auto-encoder – Experimental Results

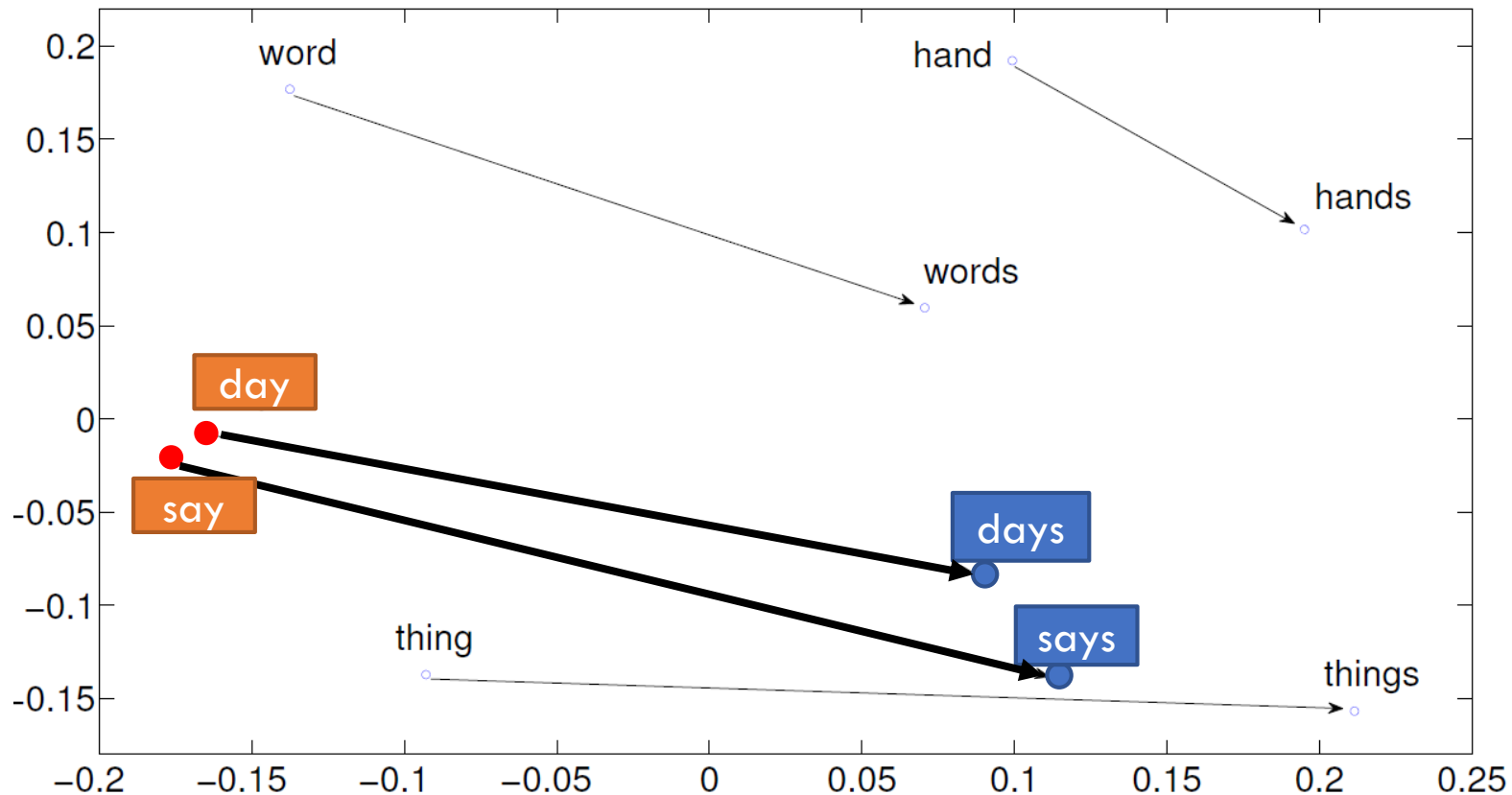


Experimental Results



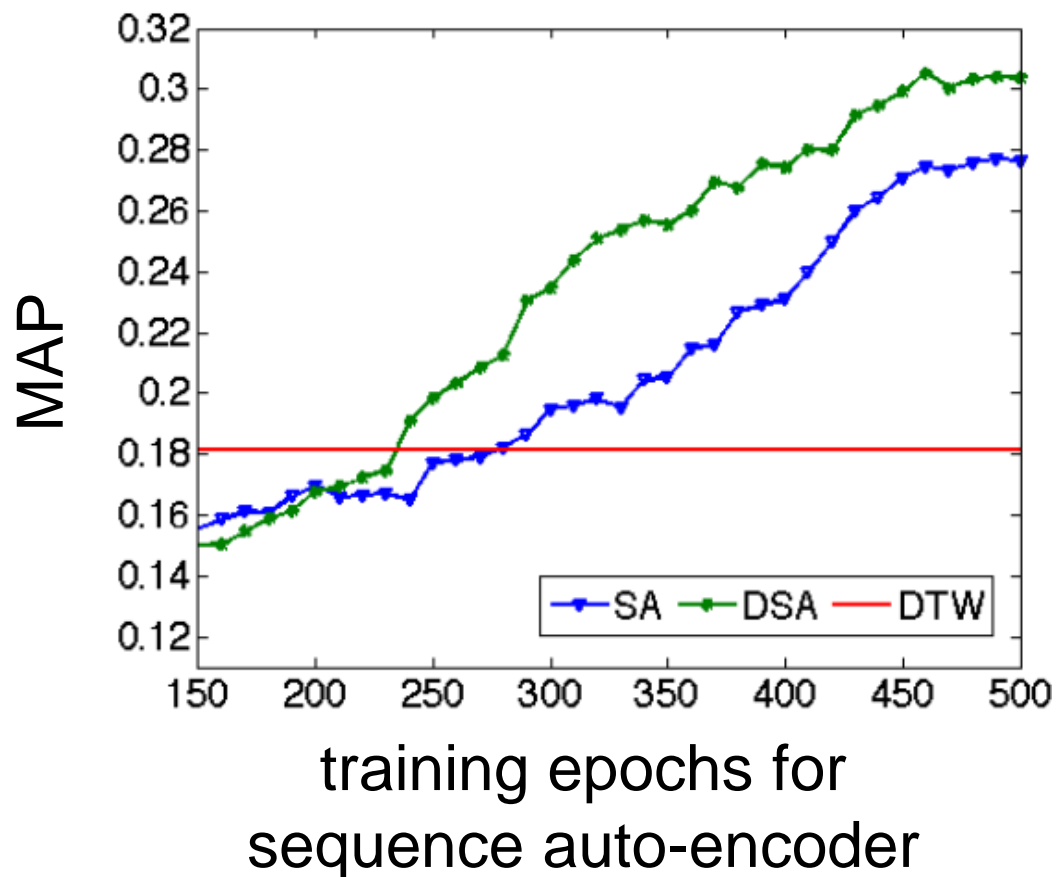
Sequence Auto-encoder – Experimental Results

- Projecting the embedding vectors to 2-D



Sequence Auto-encoder – Experimental Results

- Audio story (LibriSpeech corpus)



SA: sequence auto-encoder

DSA: de-noising sequence auto-encoder

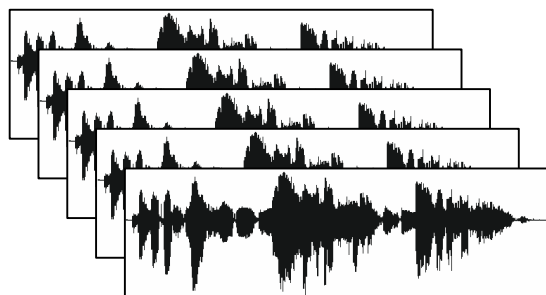
Input: clean speech
+ noise

output: clean speech

New Direction 3-3:
No Speech Recognition!
Unsupervised ASR

Conventional ASR

unknown speech signal



ASR

... Hello World ...

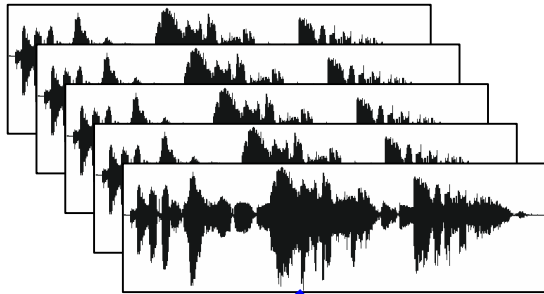
A huge annotated
corpus

Knowledge about the
language (Phone set, Lexicon,
Language Model)



Unsupervised ASR

unknown speech signal



ASR

Acoustic Tokens

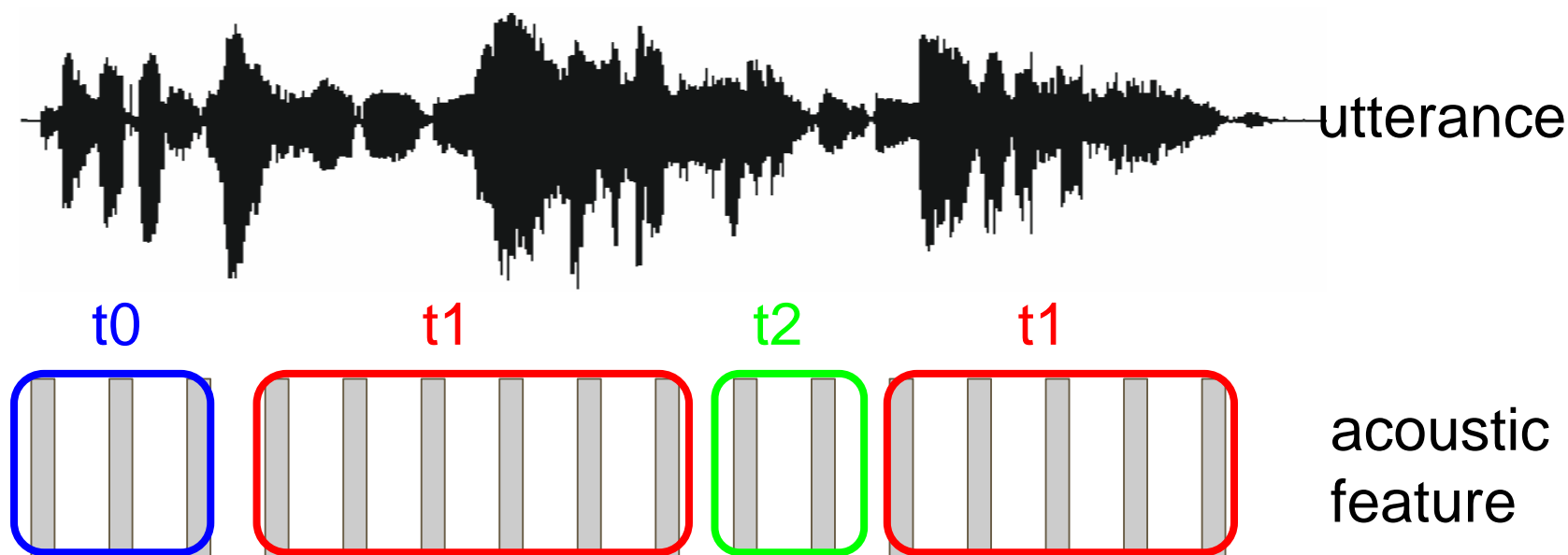
$t_0t_1t_2, t_1t_3, t_2t_3,$
 $t_2t_1t_3t_3t_2 \dots$

Used in Query by example
Spoken Term Detection

Unsupervised ASR:

Learn the models for a set of acoustic patterns (tokens)
directly from the corpus (target spoken archive)

Unsupervised ASR - Acoustic Token



acoustic tokens: chunks of acoustically similar feature vectors with token ids

[Zhang & Glass, ASRU 09]
[Huijbregts, ICASSP 11]
[Chan & Lee, Interspeech 11]

Unsupervised ASR

- Overall Framework

$$\omega_0 = \text{initialization}(X)$$

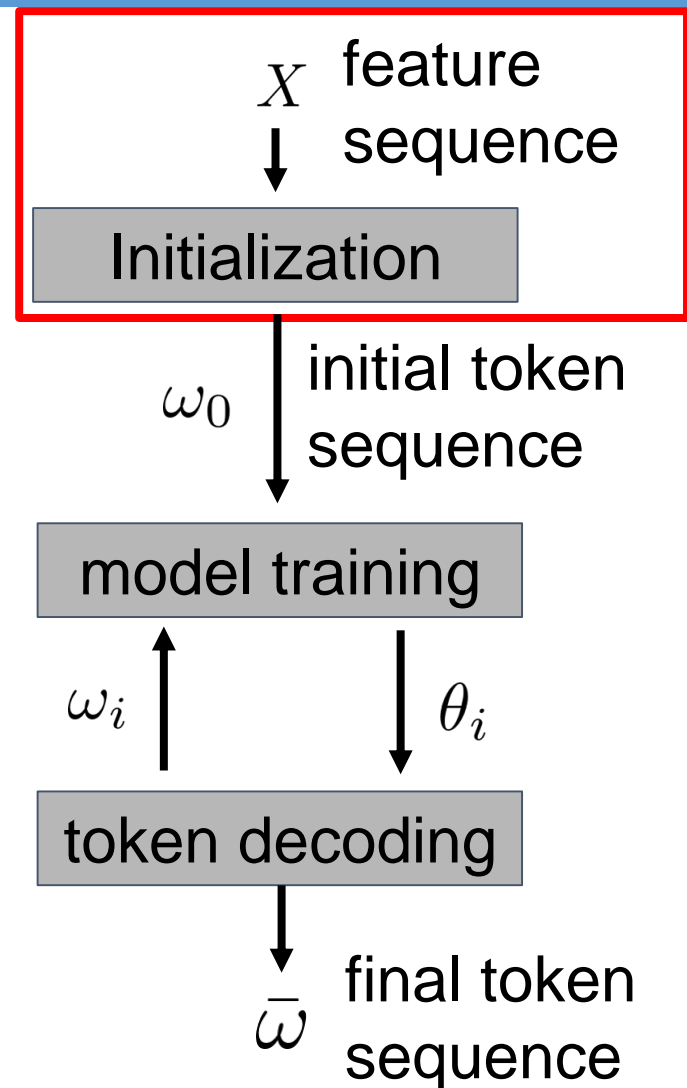
**simple segmentation
and clustering**

$$\theta_i = \arg \max_{\theta} P(X | \omega_{i-1}, \theta)$$

$$\omega_i = \arg \max_{\omega} P(X | \omega, \theta_{i-1})$$

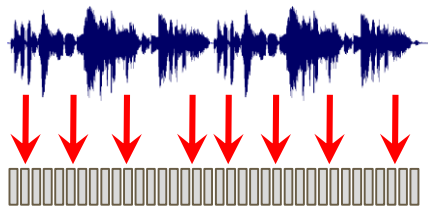
X : feature sequence for the whole corpus

ω : token sequences for X

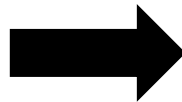


Unsupervised ASR

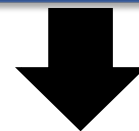
- Initialization



Extract acoustic features for every utterance



Grouping consecutive acoustically similar feature vectors into segments



Extract mean of each segment and perform K-means clustering on the entire archive



Get Token ID



Unsupervised ASR

- Overall Framework

$$\omega_0 = \text{initialization}(X)$$

**simple segmentation
and clustering**

$$\theta_i = \arg \max_{\theta} P(X | \omega_{i-1}, \theta)$$

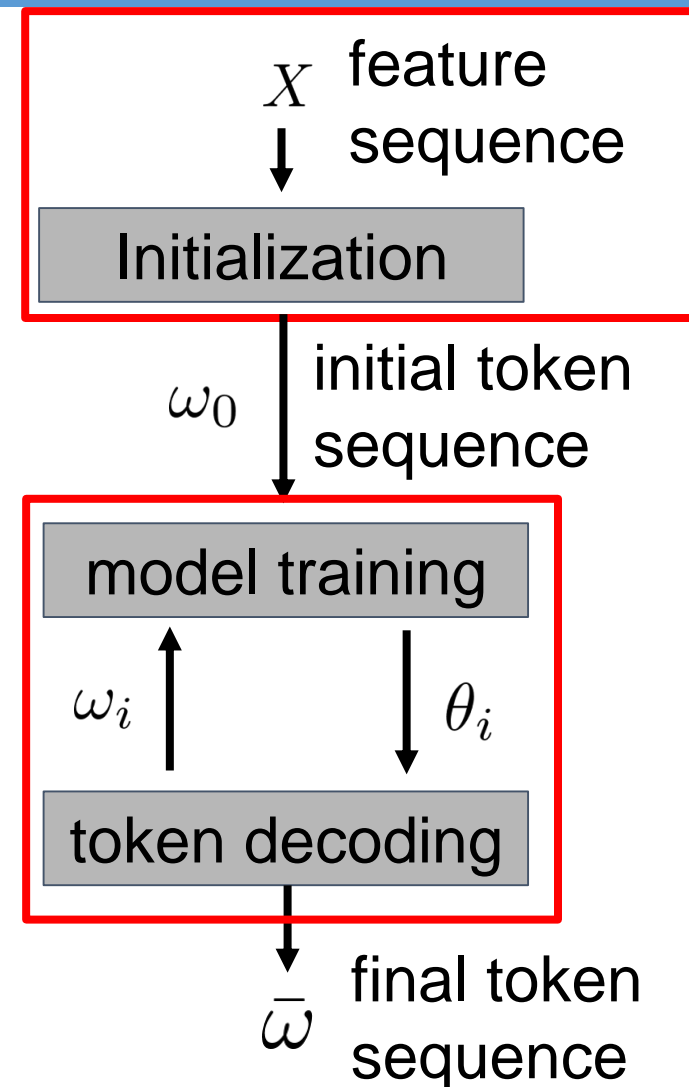
$$\omega_i = \arg \max_{\omega} P(X | \omega, \theta_{i-1})$$

X : feature sequence for the whole corpus

θ : Model (e.g. HMM) parameters

ω : token sequences for X

i : training iteration



Unsupervised ASR

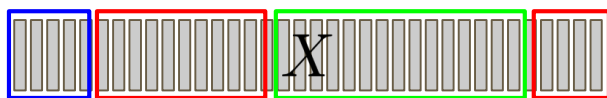
- Overall Framework

optimize HMM parameters using Baum–Welch algorithm on token sequence ω_{i-1} to get new models θ_i

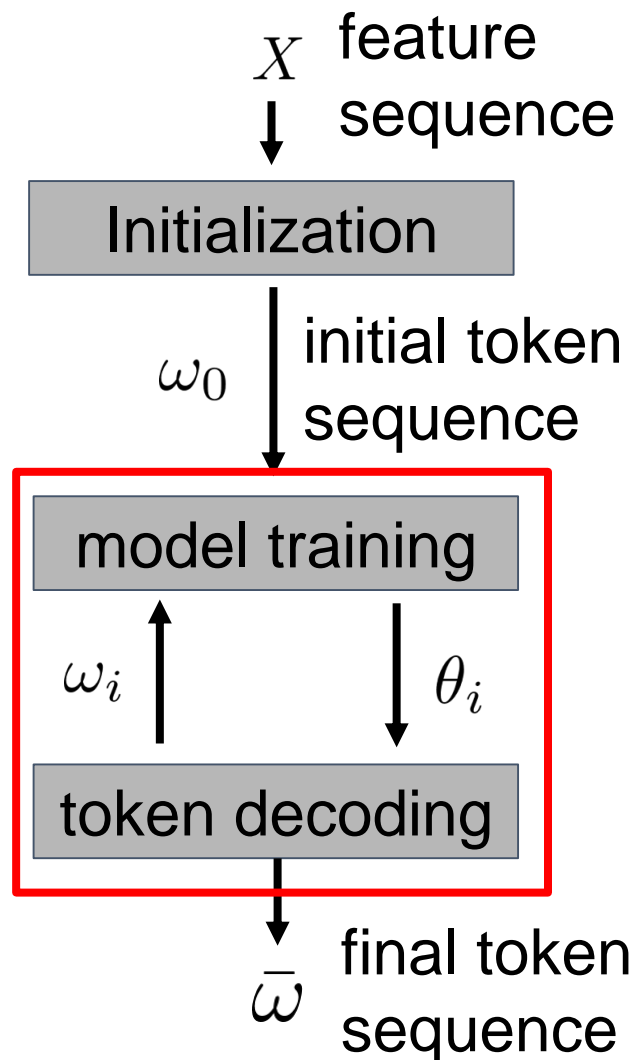


$$\theta_i = \arg \max_{\theta} P(X | \omega_{i-1}, \theta)$$

decode acoustic features into a new token sequence ω_i using Viterbi decoding

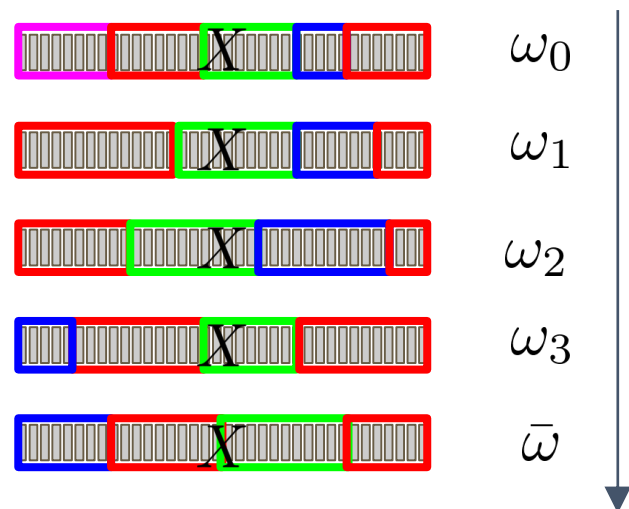


$$\omega_i = \arg \max_{\omega} P(X | \omega, \theta_{i-1})$$

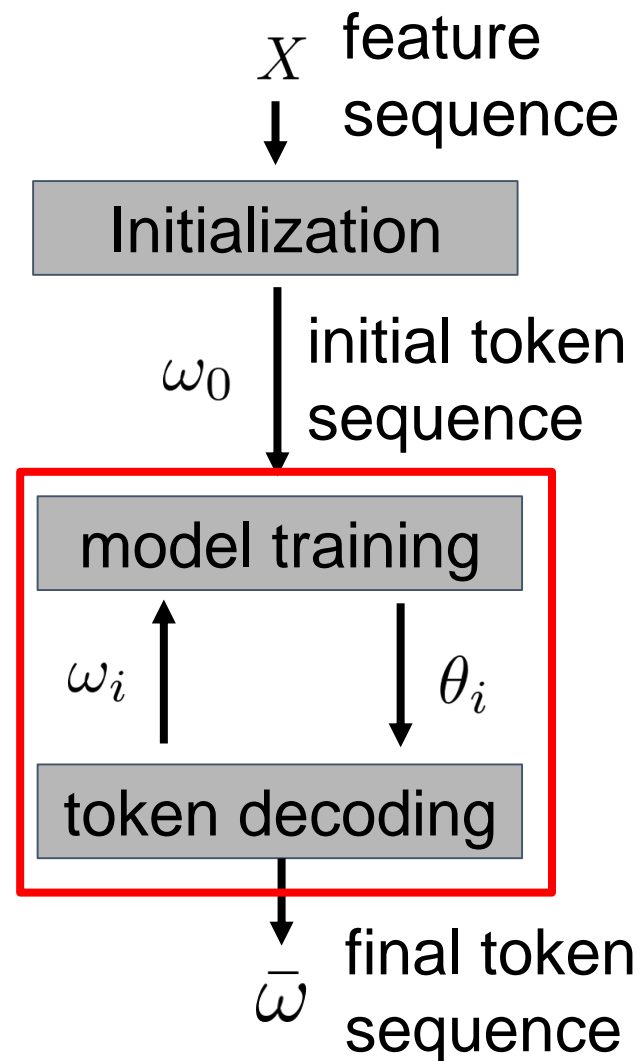


Unsupervised ASR

- Overall Framework

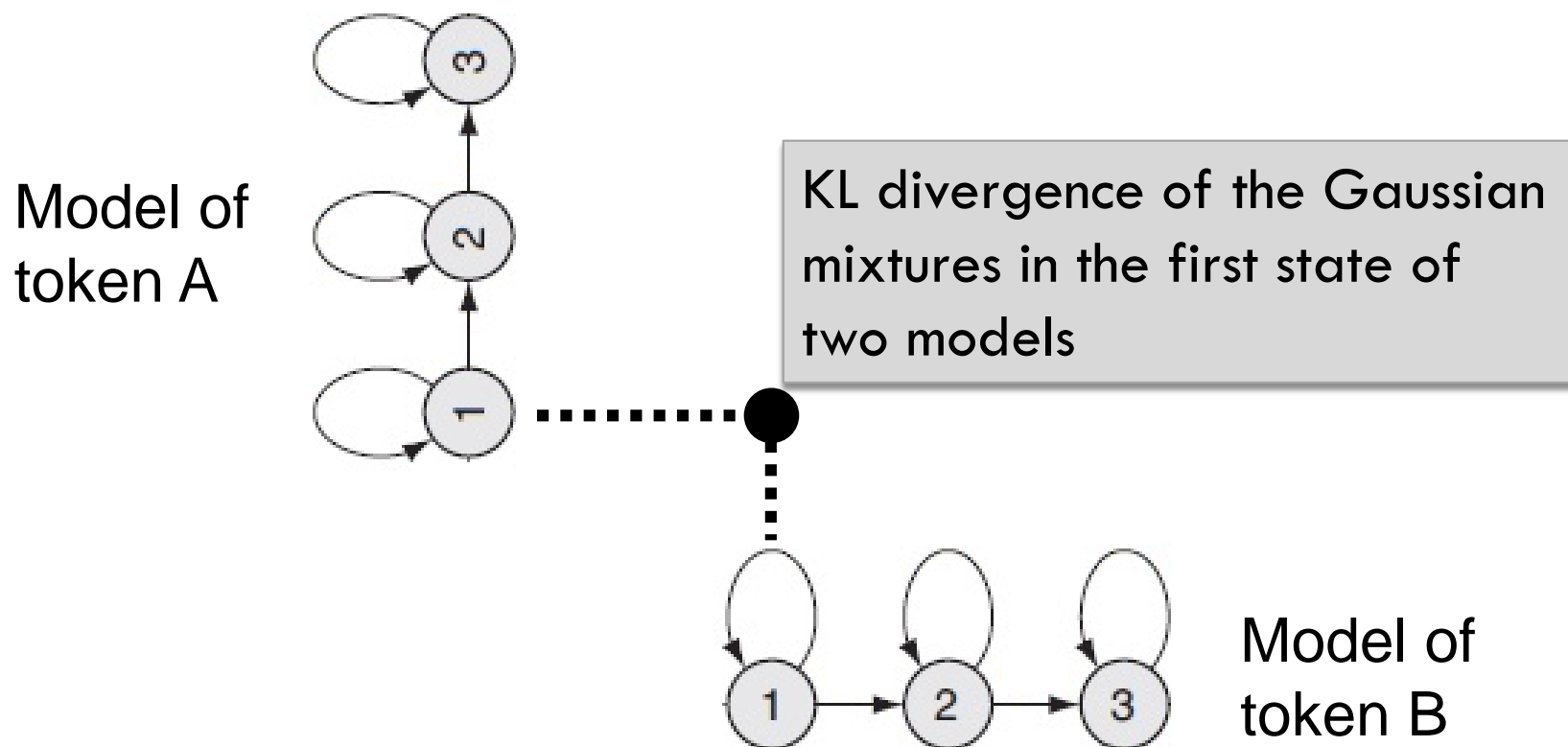


iterate until the token sequences (including token boundaries) converge



Acoustic Token in Query by Example Spoken Term Detection

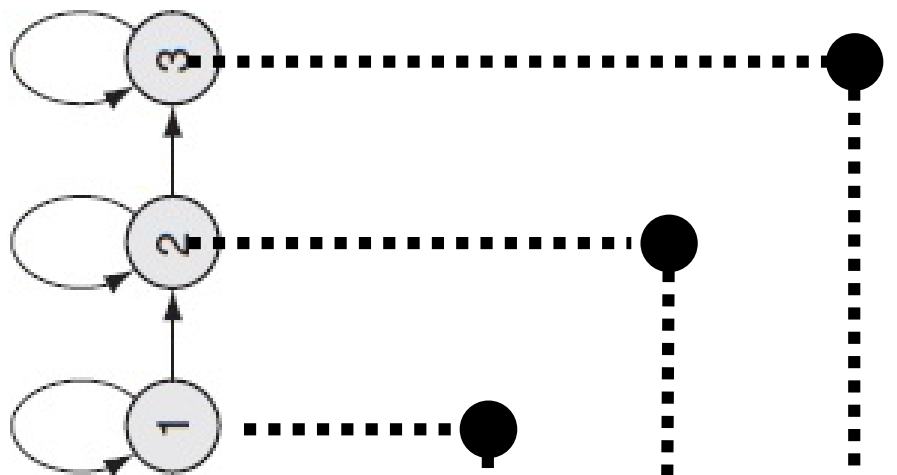
- Compute the similarity between the models of two tokens



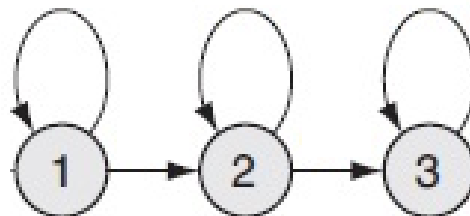
Acoustic Token in Query by Example Spoken Term Detection

- Compute the similarity between the models of two tokens

Model of token A



Sum of the KL divergence over the states of the two token models

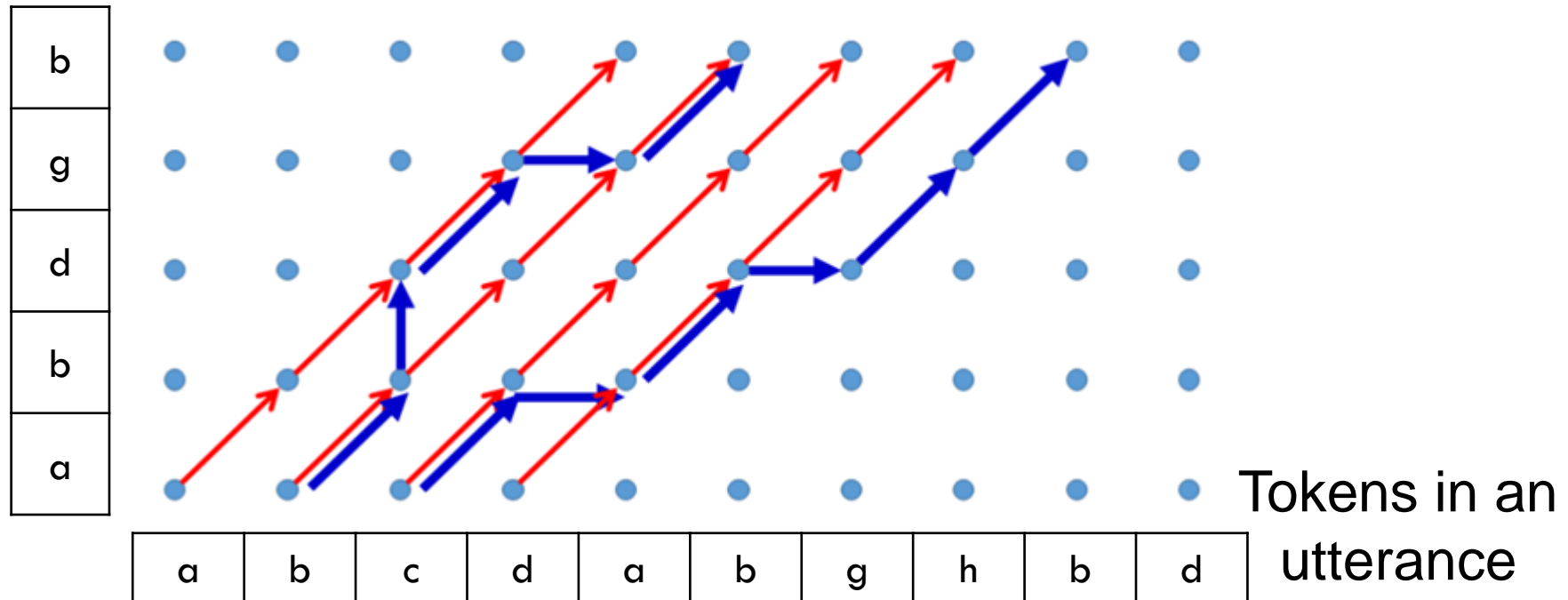


Model of token B

Token-based DTW

Tokens
in query

subsequence matching Token-based DTW



- Signal-level DTW is more sensitive to signal variation (e.g. same phoneme across different speakers), while token models are able to cover better the distribution of signal variation
- Much lower on-line computation load

Multi-granularity Space for Acoustic Tokens

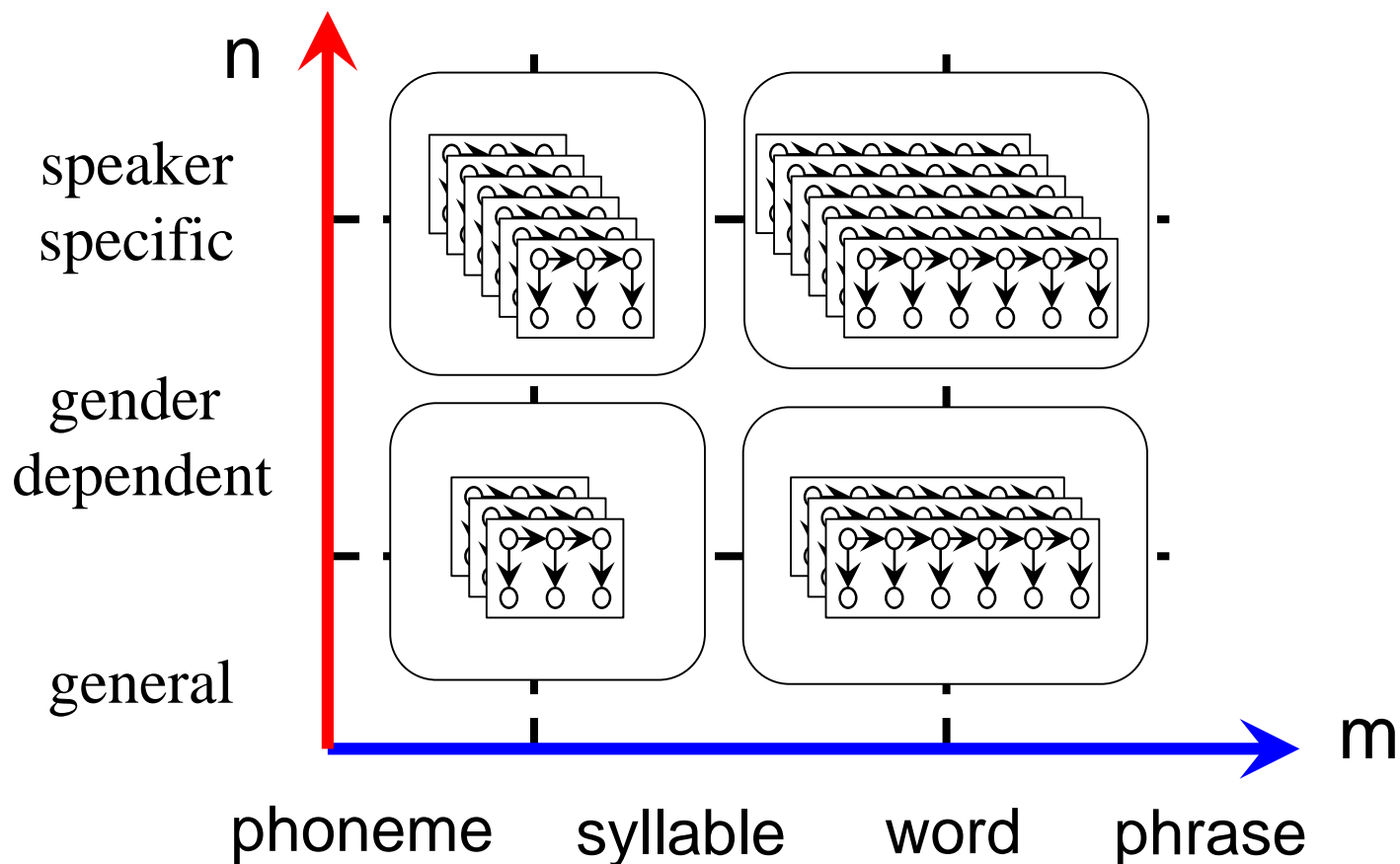
- Unknown hyperparameters for the token models
 - Number of HMM states per token (m): token length
 - Number of distinct tokens (n)
- Multiple layers of intrinsic representations of speech

Multi-granularity Space for Acoustic Tokens

- From short to long (Temporal Granularity)
 - phoneme Number of states per HMM (m)
 - syllable
 - word
 - phrase
- From coarse units to fine units (Phonetic Granularity)
 - general phoneme set Number of distinct HMMs (n)
 - gender dependent phoneme set
 - speaker specific phoneme set

Multi-granularity Space for Acoustic Tokens

Training multiple sets of HMMs for with different granularity
[Chung & Lee, ICASSP 14]



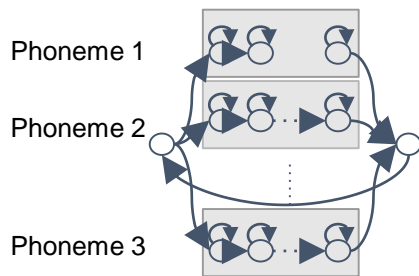
Multi-granularity Space for Acoustic Tokens

- Token-based DTW using tokens with different granularity (m,n) averaged gave much better performance
- One example
 - ▣ Frame-level DTW: $MAP = 10\%$
 - ▣ Using only the token set with the best performance: $MAP = 11\%$
 - ▣ Using 20 sets of tokens (number of states per HMM $m = 3,5,7,9,11$, number of distinct HMMs $n=50,100,200,300$): $MAP = 26\%$

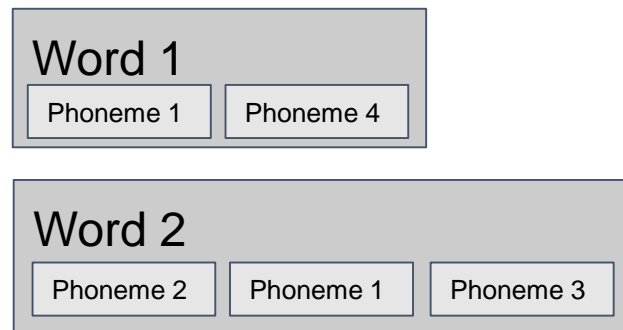
Hierarchical Paradigm

- Typical ASR:
 - ▣ Acoustic Model: models for the phonemes
 - ▣ Lexicon: the pronunciation of every word as a phoneme sequence
 - ▣ Language Model: the transition between words

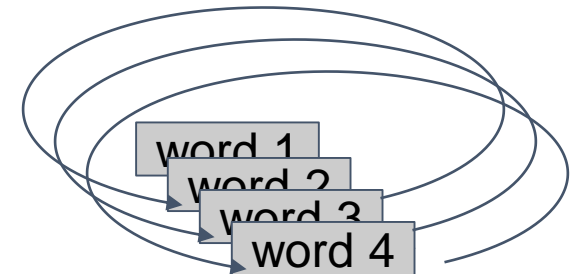
Acoustic Model



Lexicon



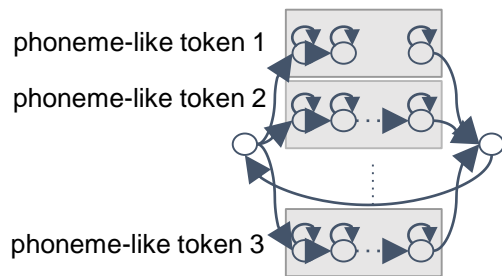
Language Model



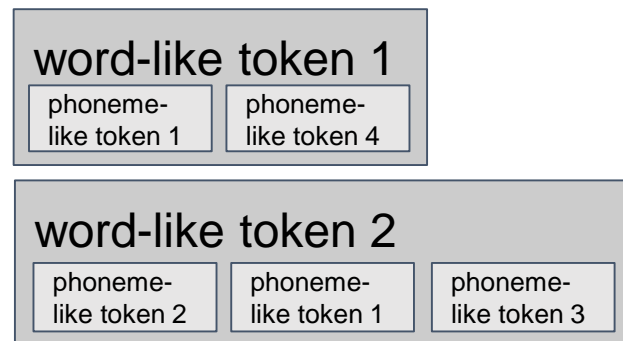
Hierarchical Paradigm

- Similarly, in unsupervised ASR:
 - ⌘ Acoustic Model: the phoneme-like token HMMs
 - ⌘ Lexicon: the pronunciation of every word-like token as a sequence of phoneme-like tokens
 - ⌘ Language Model: the transition between word-like tokens

Acoustic Model



Lexicon



Language Model



Hierarchical Paradigm

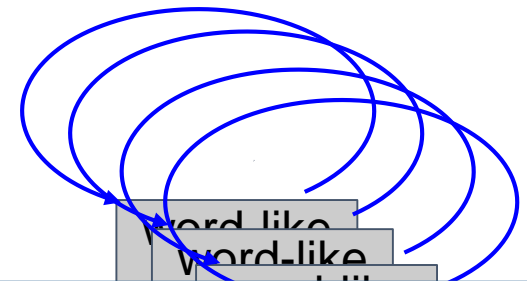
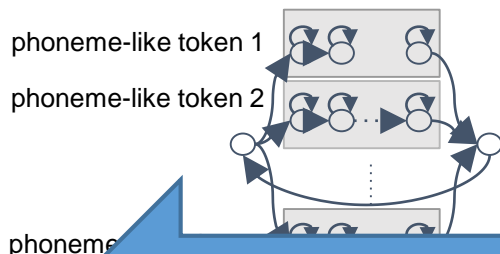
- Similarly, in unsupervised ASR:
 - ⌘ Acoustic Model: the phoneme-like token HMMs
 - ⌘ Lexicon: the pronunciation of every word-like token as a sequence of phoneme-like tokens
 - ⌘ Language Model: the transition between word-like tokens

Bottom-up Construction

Acoustic Model

Lexicon

Language Model



Top Down Constraint

Bottom Up Construction

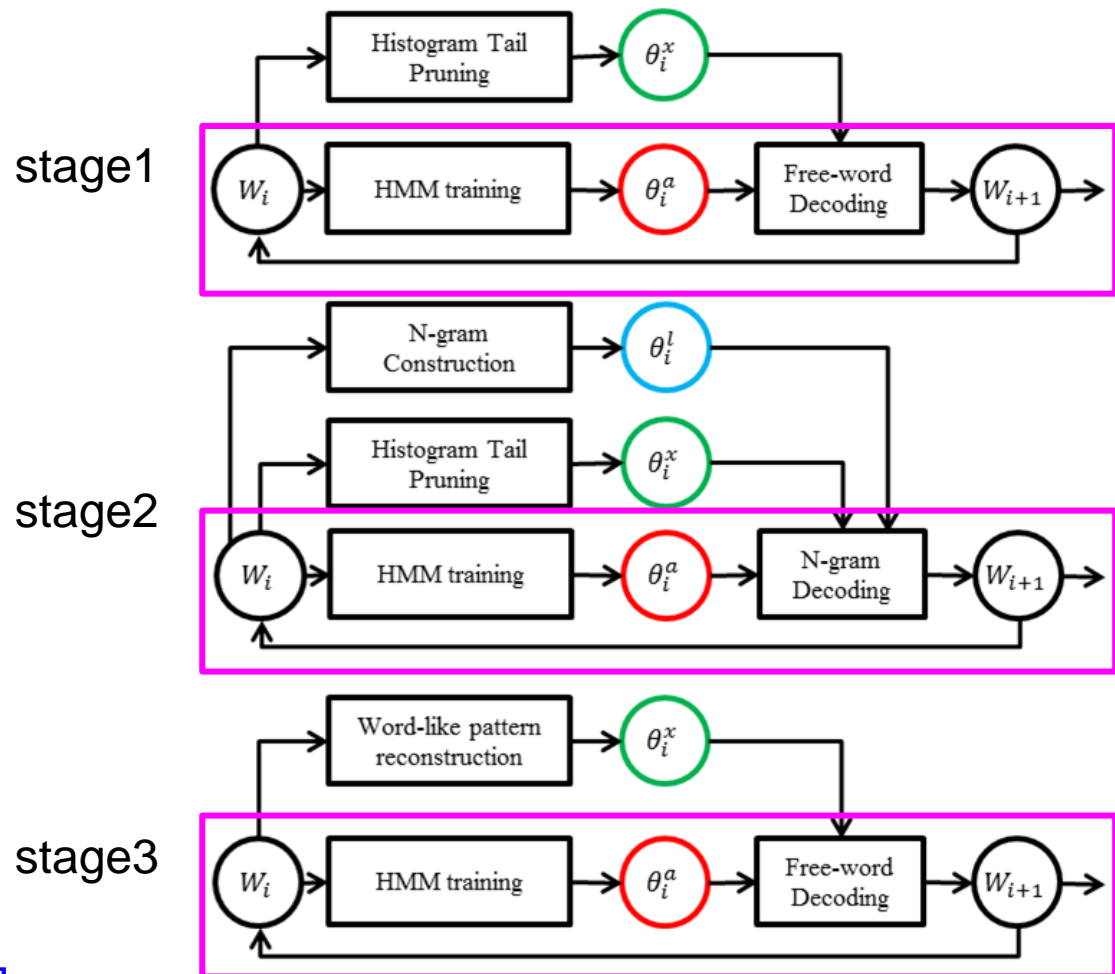
3 stages during training
focus on different constraints:

stage1: **Acoustic Model** θ_i^a

stage2: **Language Model** θ_i^l

stage3: **Lexicon** θ_i^x

this part alone would be the
HMM training we described
earlier



Hierarchical Paradigm

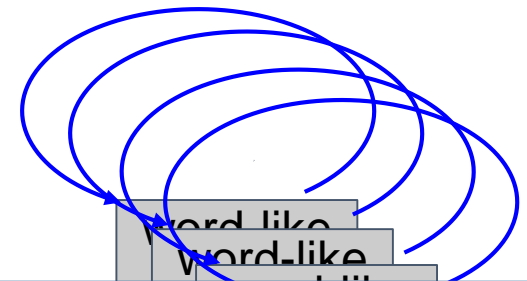
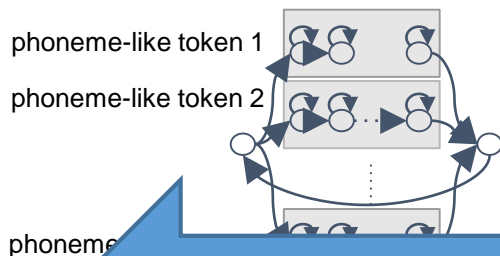
- Similarly, in unsupervised ASR:
 - ⌘ Acoustic Model: the phoneme-like token HMMs
 - ⌘ Lexicon: the pronunciation of every word-like token as a sequence of phoneme-like tokens
 - ⌘ Language Model: the transition between word-like tokens

Bottom-up Construction

Acoustic Model

Lexicon

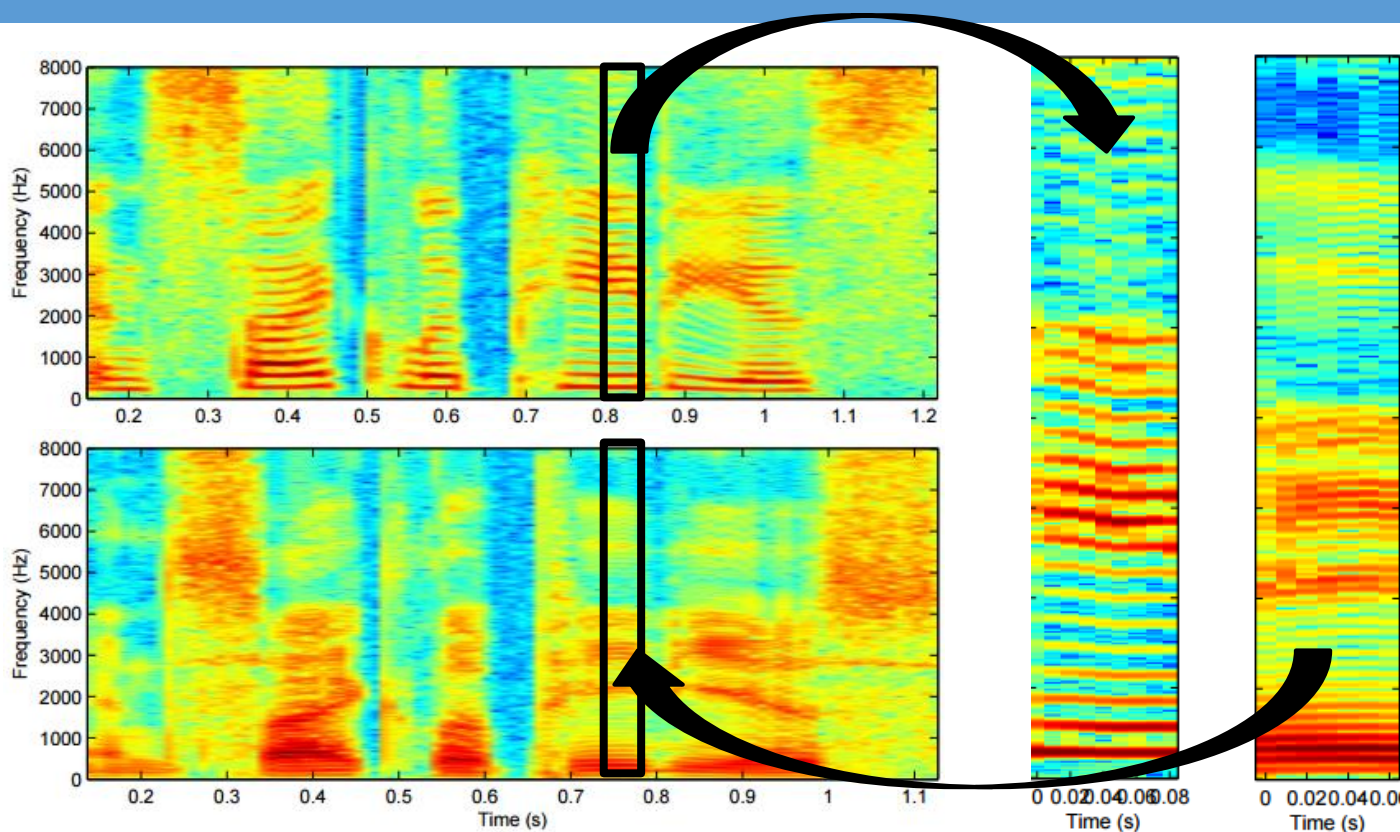
Language Model



Top Down Constraint

Top-down Constraints

This figure is from Aren Jansen's ICASSP paper.
[Jansen, ICASSP 13]



- Signals of the same phoneme may be very different on phoneme level, but the global structures of signals of the same word are very often very similar on word level
- Global structures help in building the hierarchical model

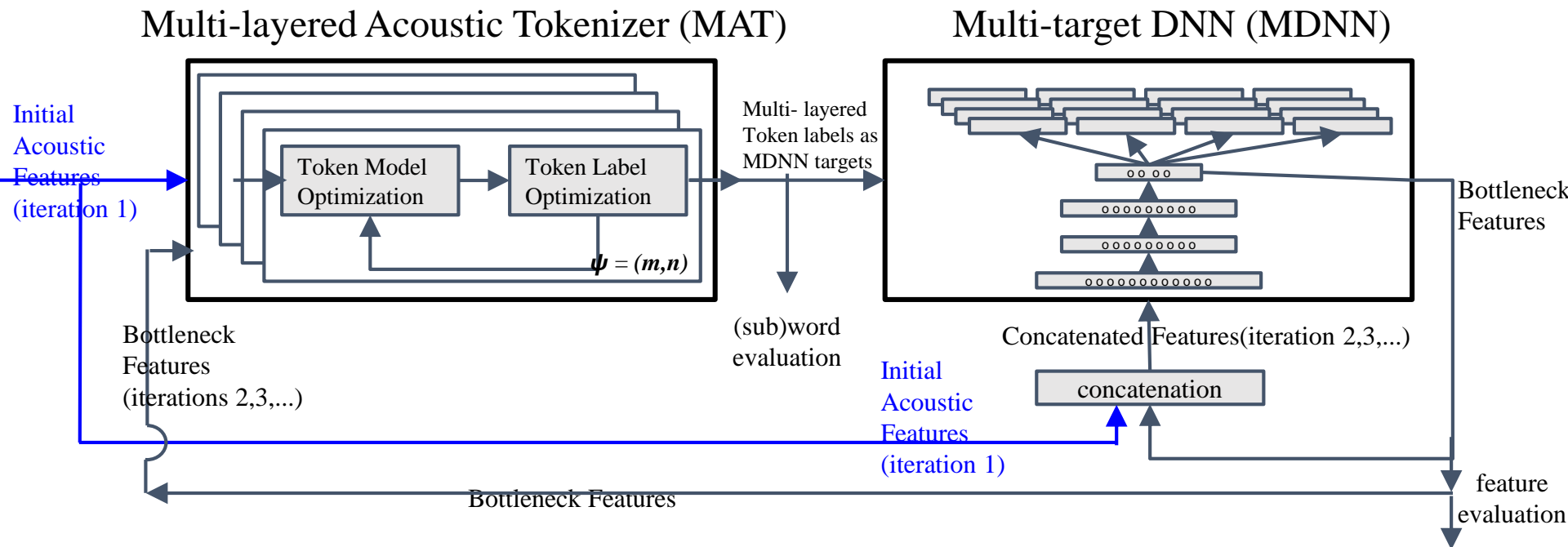
Multi-layered Acoustic Tokenizing Deep Neural Networks (MAT-DNN)

[Chung & Lee, ASRU 15]

- Jointly learn high quality frame-level features (much better than MFCCs) and acoustic tokens in an unsupervised way
- Unsupervised training of multi-target DNN using unsupervised token labels as training target

In the **first iteration**, we use MFCC as the initial features

In the other iterations, we concatenate the bottleneck features with the MFCC



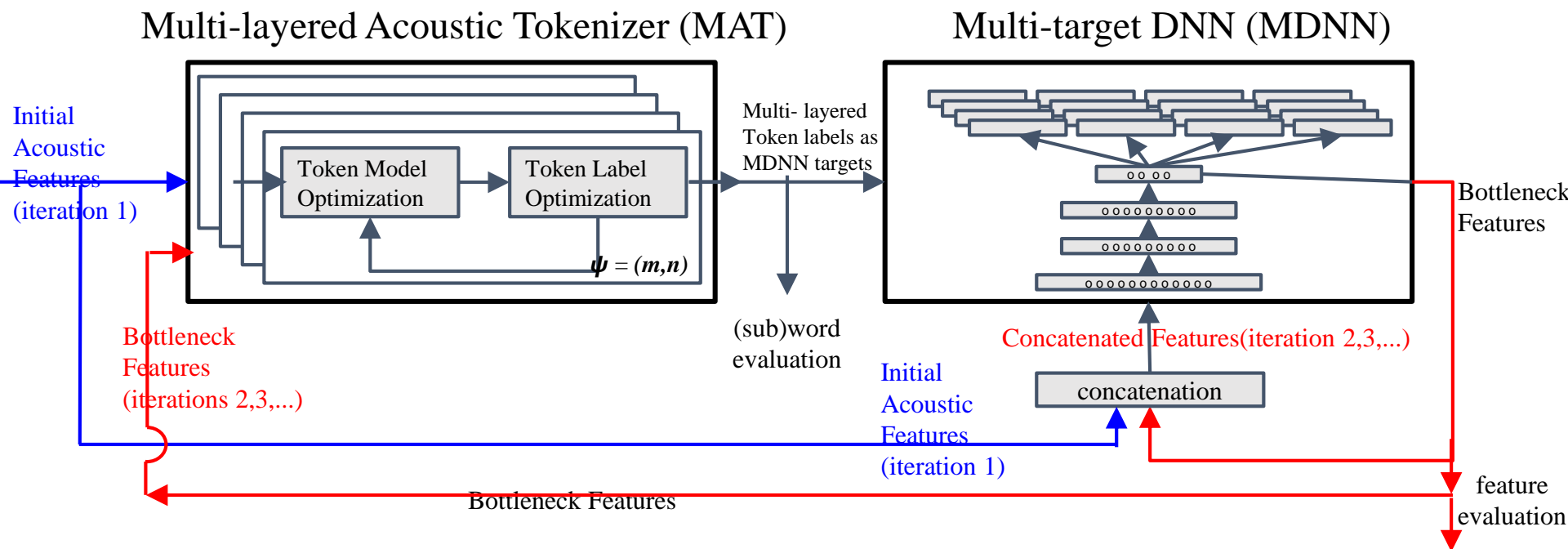
Multi-layered Acoustic Tokenizing Deep Neural Networks (MAT-DNN)

[Chung & Lee, ASRU 15]

- Jointly learn high quality frame-level features (much better than MFCCs) and acoustic tokens in an unsupervised way
- Unsupervised training of multi-target DNN using unsupervised token labels as training target

In the **first iteration**, we use MFCC as the initial features

In the **other iterations**, we concatenate the bottleneck features with the MFCC



Multi-layered Acoustic Tokenizing Deep Neural Networks (MAT-DNN)

[Chung & Lee, ASRU 13]

- Experimental Results
 - ▣ Query by Example Spoken Term Detection on Tsonga

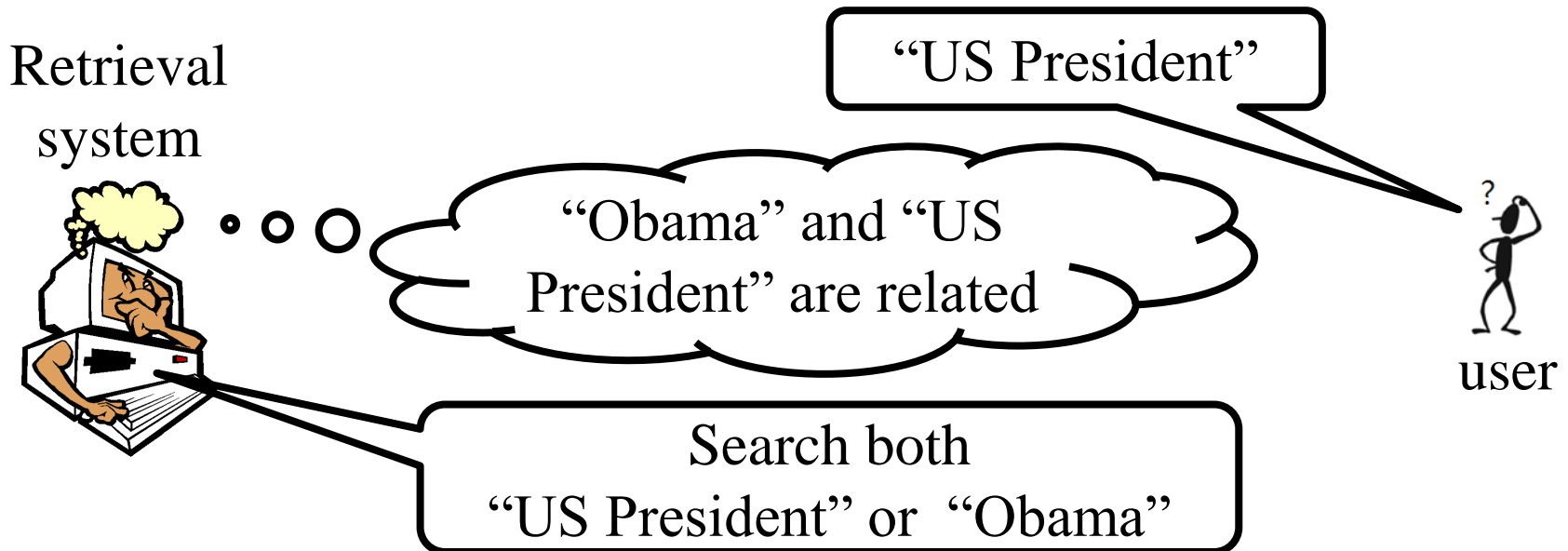
Approach		MAP
Frame-based DTW	MFCC	9.0
	New Feature	28.7
Token-based DTW	New Tokens	26.2

New Direction 4:
Special Semantic Retrieval Techniques
for Spoken Content



Semantic Retrieval

- User expects semantic retrieval of spoken content.
 - ▣ User enters “US President”, system also finds “Obama”
- Widely studies on text retrieval
 - ▣ Take *query expansion* as example

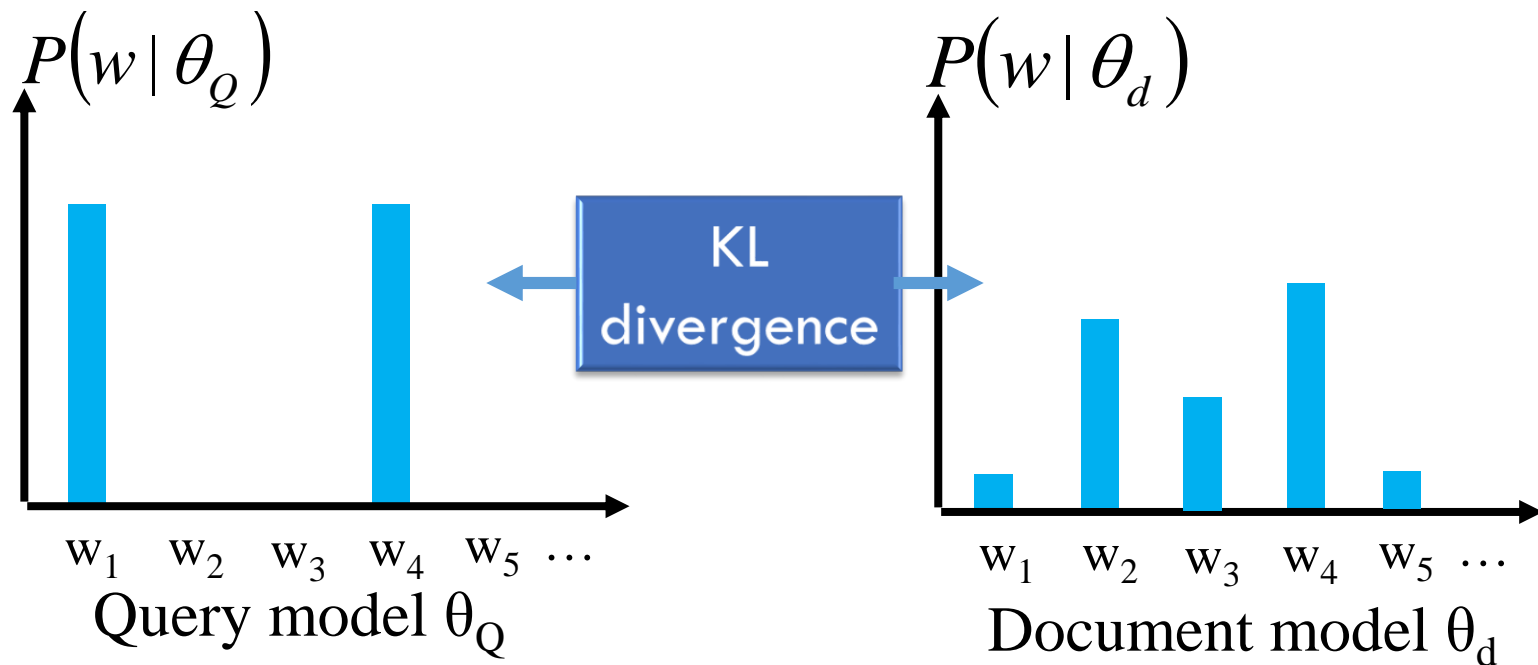


Semantic Retrieval

- User expects semantic retrieval of spoken content.
 - ▣ User enters “US President”, system also finds “Obama”
- Widely studies on text retrieval
 - ▣ Take *query expansion* as example
- The semantic retrieval techniques developed for text can be directly applied on spoken content
- Query/document expansion based on language modeling retrieval approach as example

Review: Language Modeling Retrieval Approach

- Both query Q and document d are represented as unigram language models θ_Q and θ_d



KL divergence between the two models can be evaluated.

Review: Language Modeling Retrieval Approach

- Given query Q , rank document d according to a relevance score function $S_{LM}(Q,d)$:

$$S_{LM}(Q, d) = -KL(\theta_Q | \theta_d)$$

- Inverse of KL divergence between query model θ_Q and document model θ_d
- The documents with document models θ_d similar to query model θ_Q are more likely to be relevant.

Review: Basic Query/Document Models in Text Retrieval

□ Query model θ_Q for text:
$$P(w | \theta_Q) = \frac{N(w, Q)}{\sum_{w'} N(w', Q)}$$

$N(w, Q)$: term frequency of word w in query Q

Normalize into probability

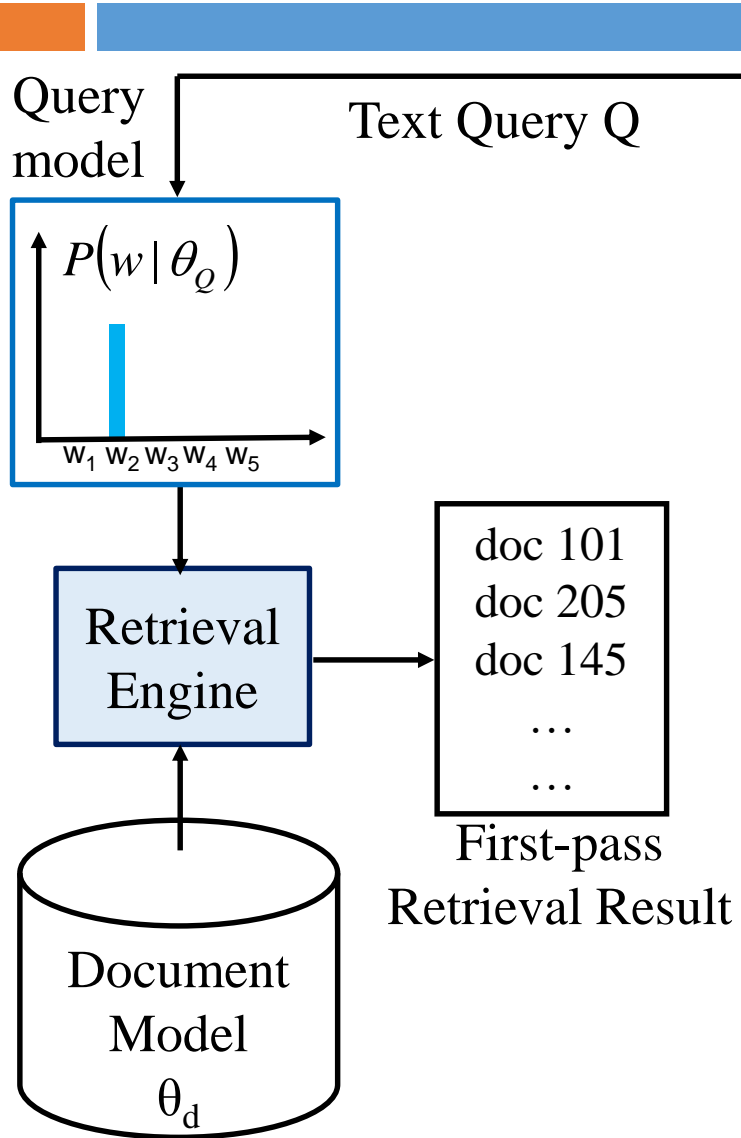
□ Document model θ_d for text:
$$P(w | \theta_d) = \frac{N(w, d)}{\sum_{w'} N(w', d)}$$

$N(w, d)$: term frequency of word w in document d

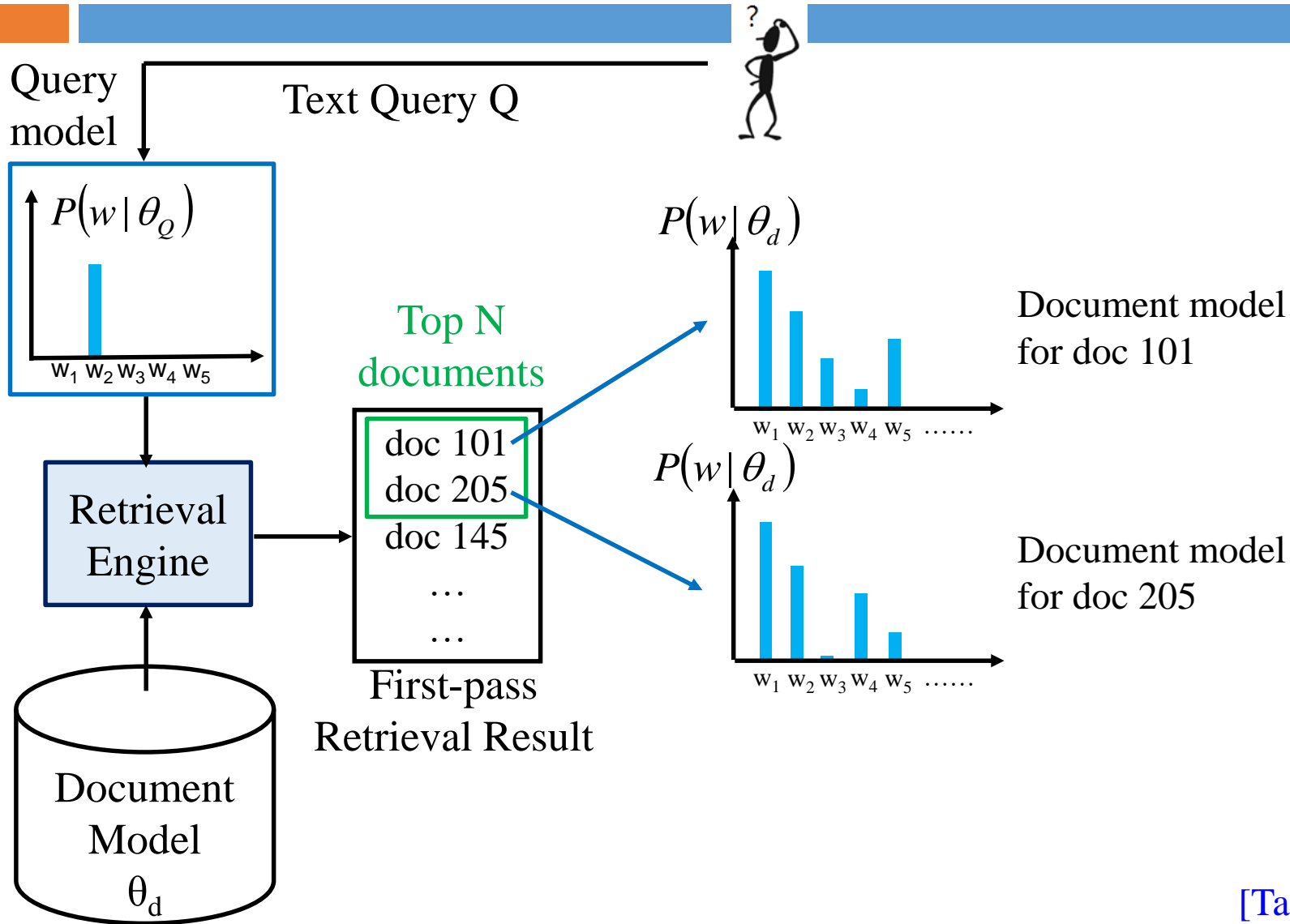
Normalize into probability

Those basic models can be enhanced by query/document expansion to handle the problem of semantic retrieval.

Review: Query Expansion

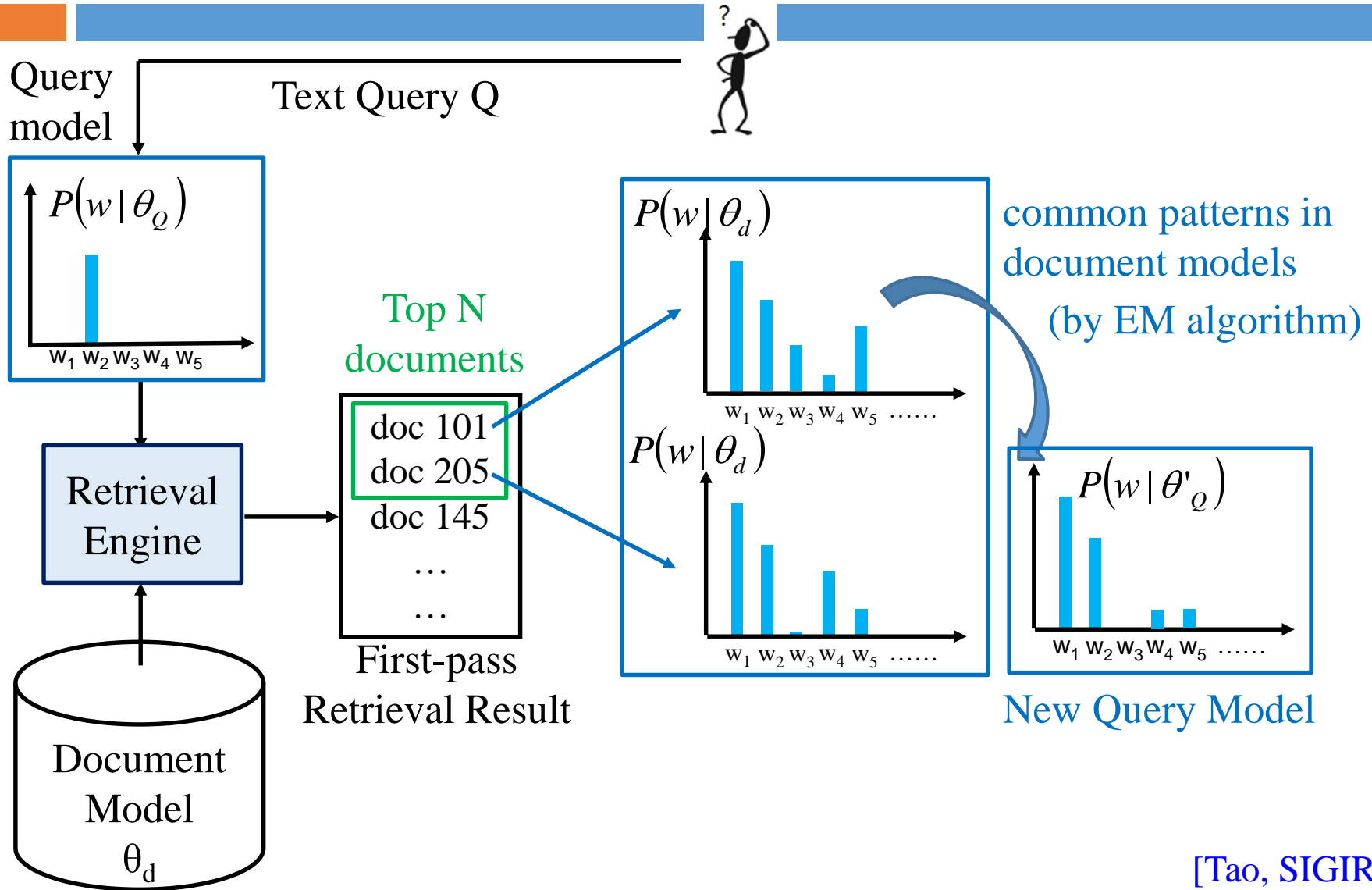


Review: Query Expansion



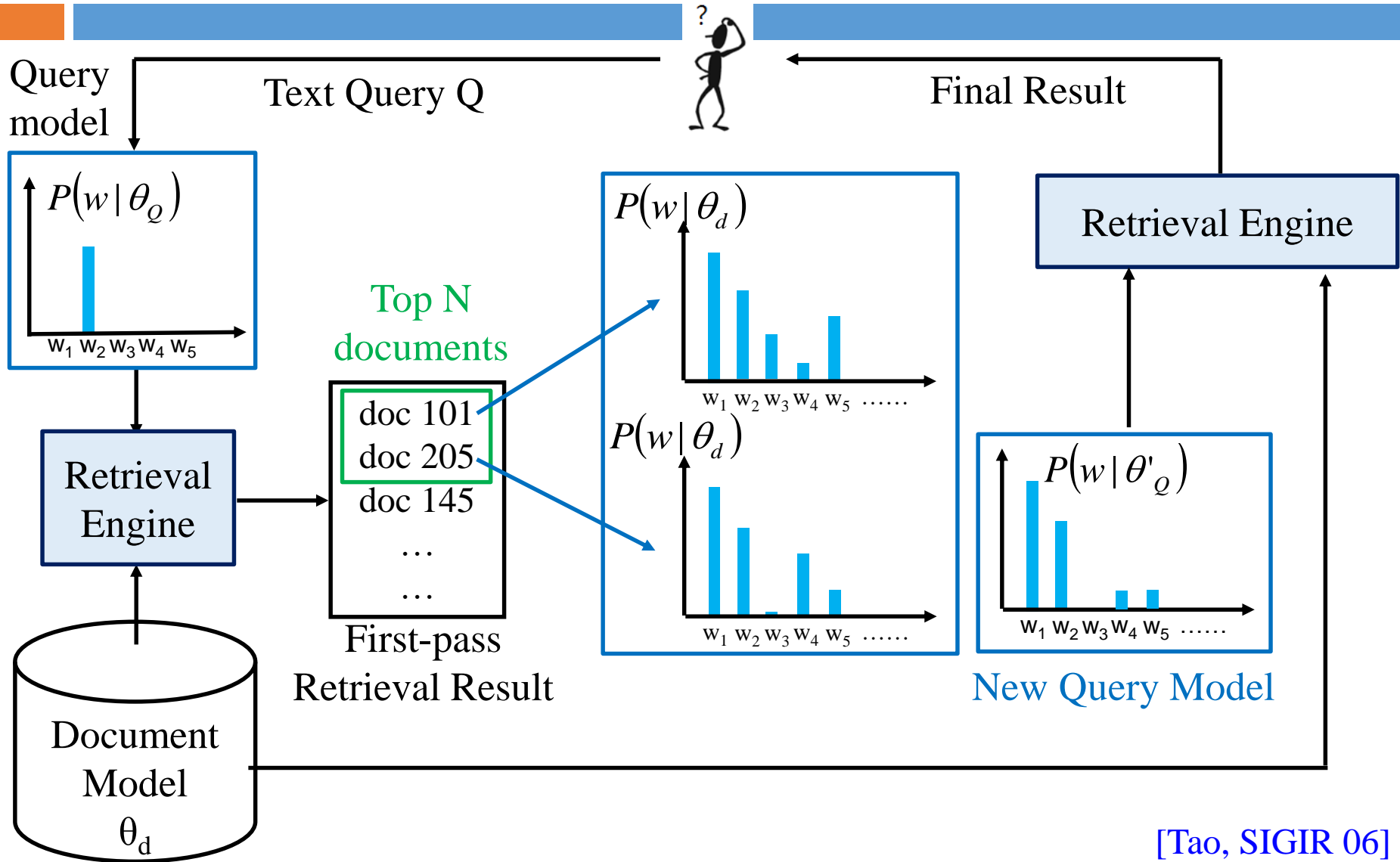
Review: Query Expansion

Parallel to
PRF

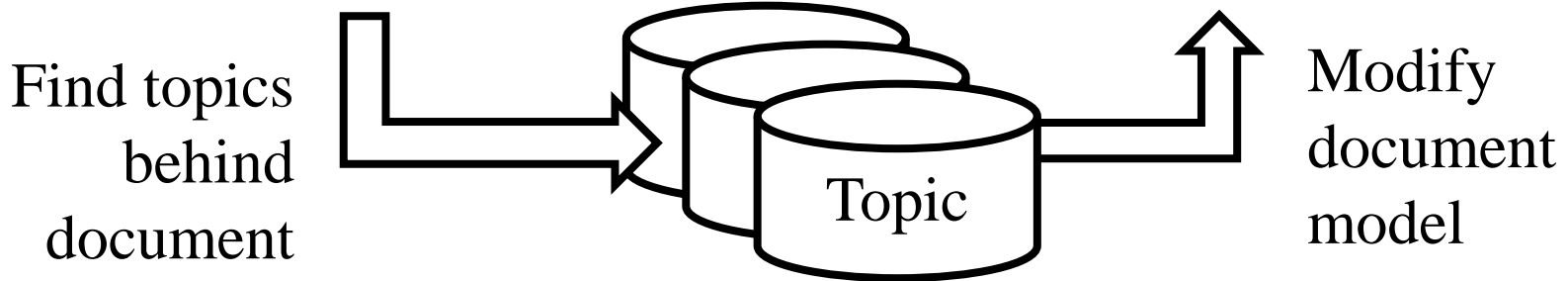
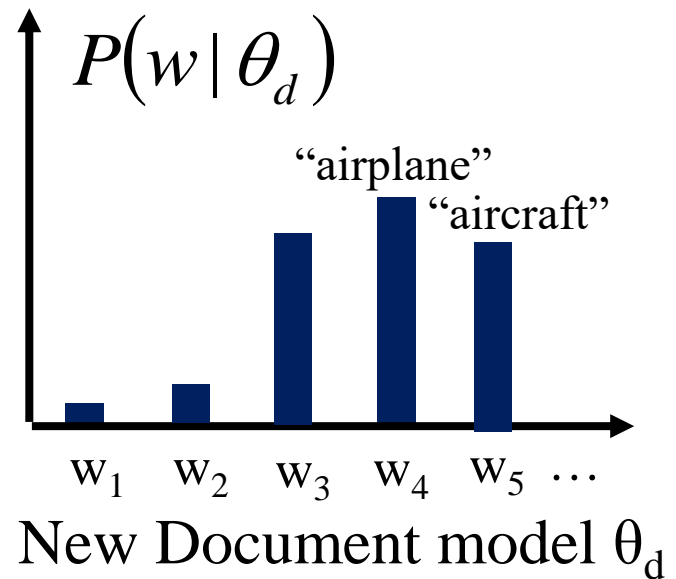
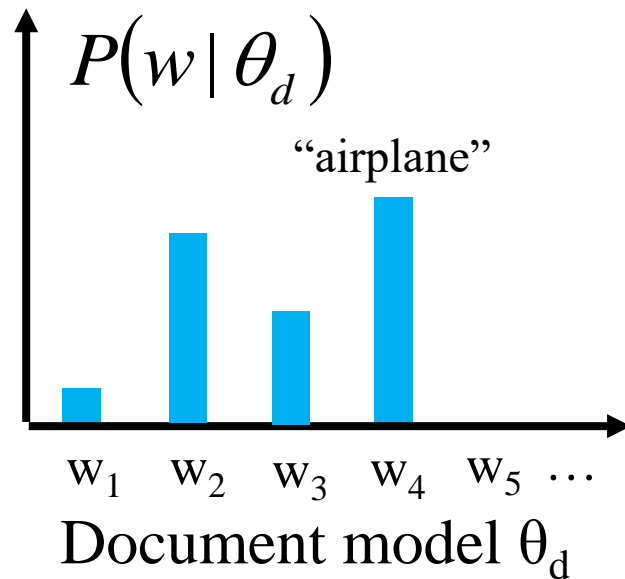


Review: Query Expansion

Parallel to PRF



Review: Document Expansion



This is realized by PLSA, LDA, etc.

Semantic Retrieval on Lattices

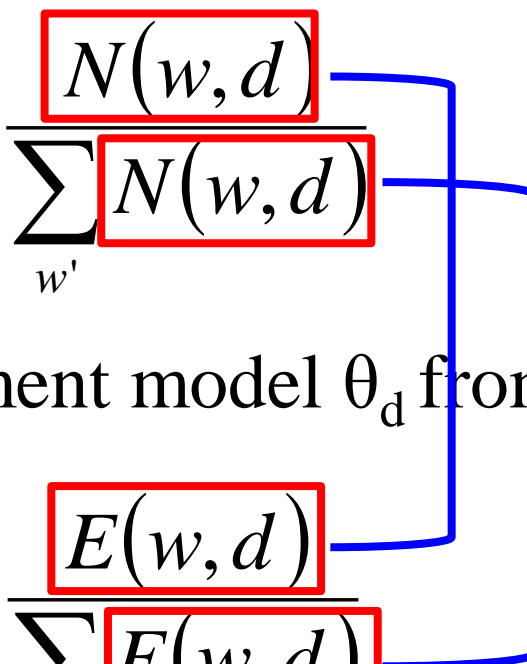
- Modify retrieval model for *lattices*:

<i>Original Retrieval Model of Text</i>		<i>For Lattices</i>
Term Frequency	➔	<i>Expected</i> Term Frequency
Document Length	➔	<i>Expected</i> Document Length
.....	

- Take the basic *language modeling retrieval approach* as example

Document Model from Lattices

- Document model θ_d for text

$$P(w | \theta_d) = \frac{N(w, d)}{\sum_{w'} N(w', d)}$$


- (Spoken) Document model θ_d from lattice

$$P(w | \theta_d) = \frac{E(w, d)}{\sum_{w'} E(w', d)}$$

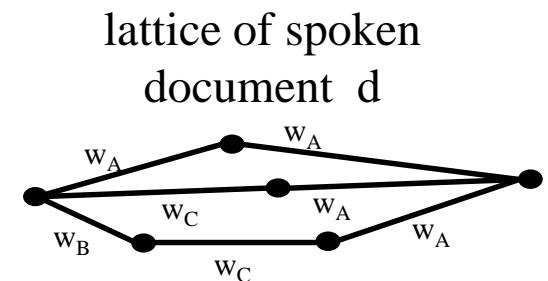
query/document expansion can be applied

Replace term frequency $N(w, d)$ with *expected term frequency $E(w, d)$ computed from lattices*

Expected Term Frequency

- **Expected** term frequency $E(w,d)$ for word w in spoken document d based on lattice

$$E(w, d) = \sum_{u \in L(d)} N(w, u) P(u | d)$$



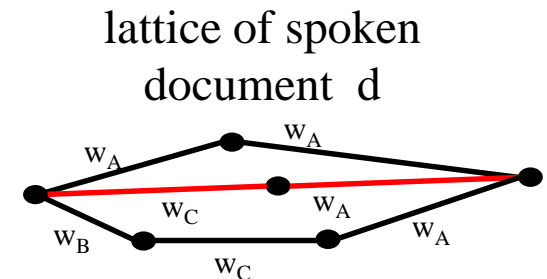
Expected Term Frequency

- **Expected** term frequency $E(w, d)$ for word w in spoken document d based on lattice

$$E(w, d) = \sum_{u \in L(d)} N(w, u) P(u | d)$$

Can we have better estimation?

- u : a word sequence in the lattice of d
- $P(u|d)$: posterior probability of word sequence u
- $N(w, u)$: the number of word w appearing in word sequence u
- $L(d)$: all the word sequences in the lattice of d



New Direction 4-1:

Special Semantic Retrieval Techniques
for Spoken Content

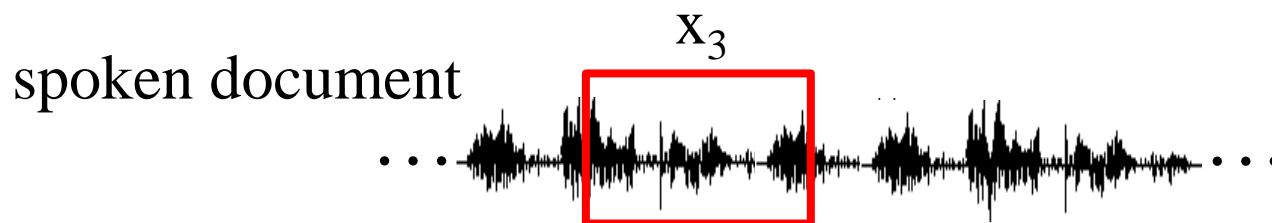
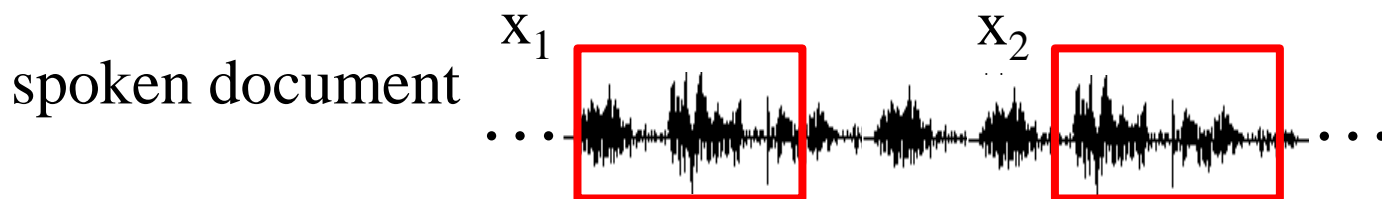
Better Estimation of Term Frequencies

Better Estimation of Term Frequencies

- Context of each term in the lattices [Tu & Lee, ICASSP 12]
 - ▣ The occurrences of a given term are usually characterized by similar context, while widely-varying contexts typically imply different terms [Schneider, Interspeech 10]
- Graph-based approach
 - ▣ Graph-based approach improved *spoken term detection*
 - ▣ It can also improve *semantic retrieval* of spoken content
 - ▣ Idea: Replace expected term frequency $E(w,d)$ with scores from graph-based approach [Lee & Lee, SLT 12] [Lee & Lee, IEEE/ACM T. ASL 14]

Graph-based Approach for Semantic Retrieval

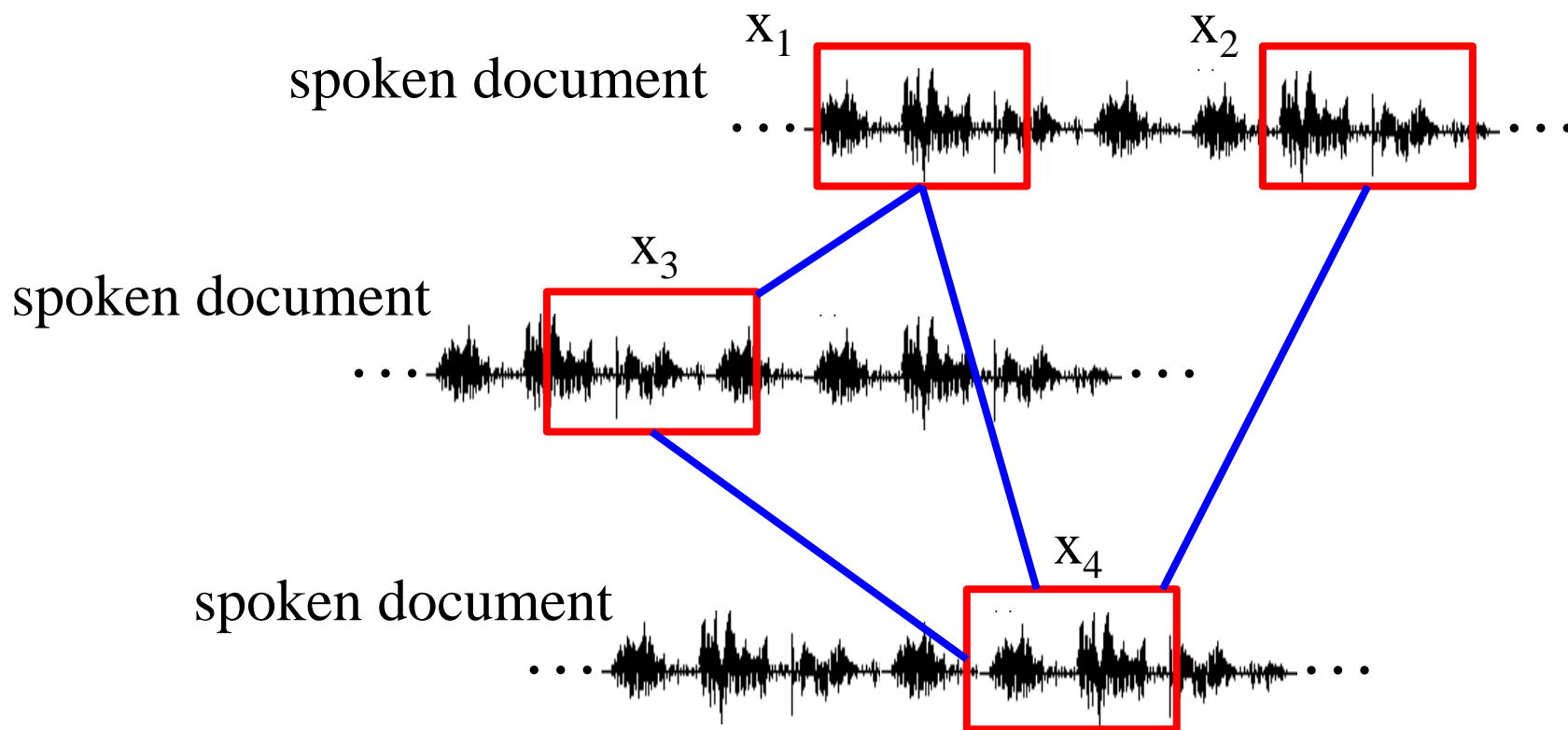
- For each word w in the lexicon



Find the occurrence regions of word w from lattices

Graph-based Approach for Semantic Retrieval

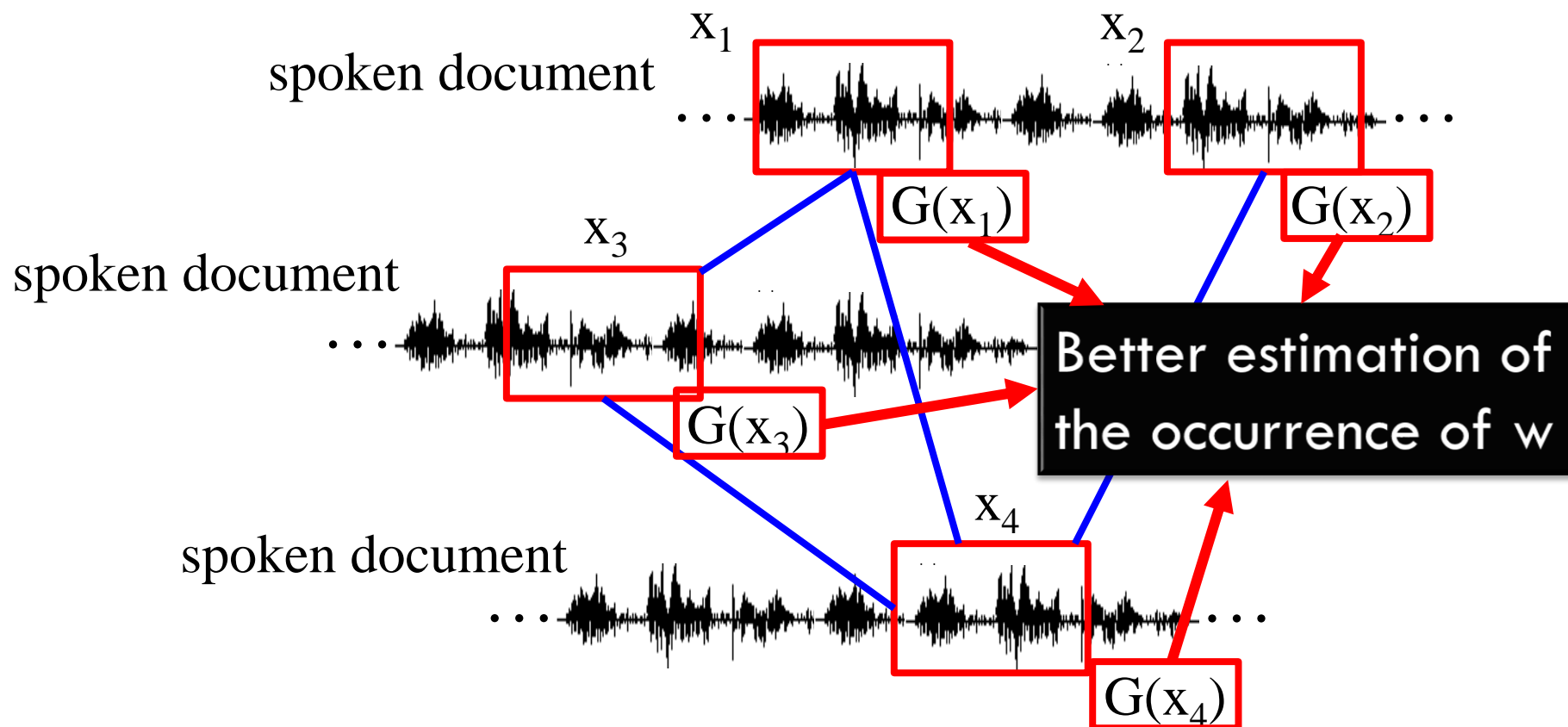
- For each word w in the lexicon



Connect the occurrence regions as a graph by similarities

Graph-based Approach for Semantic Retrieval

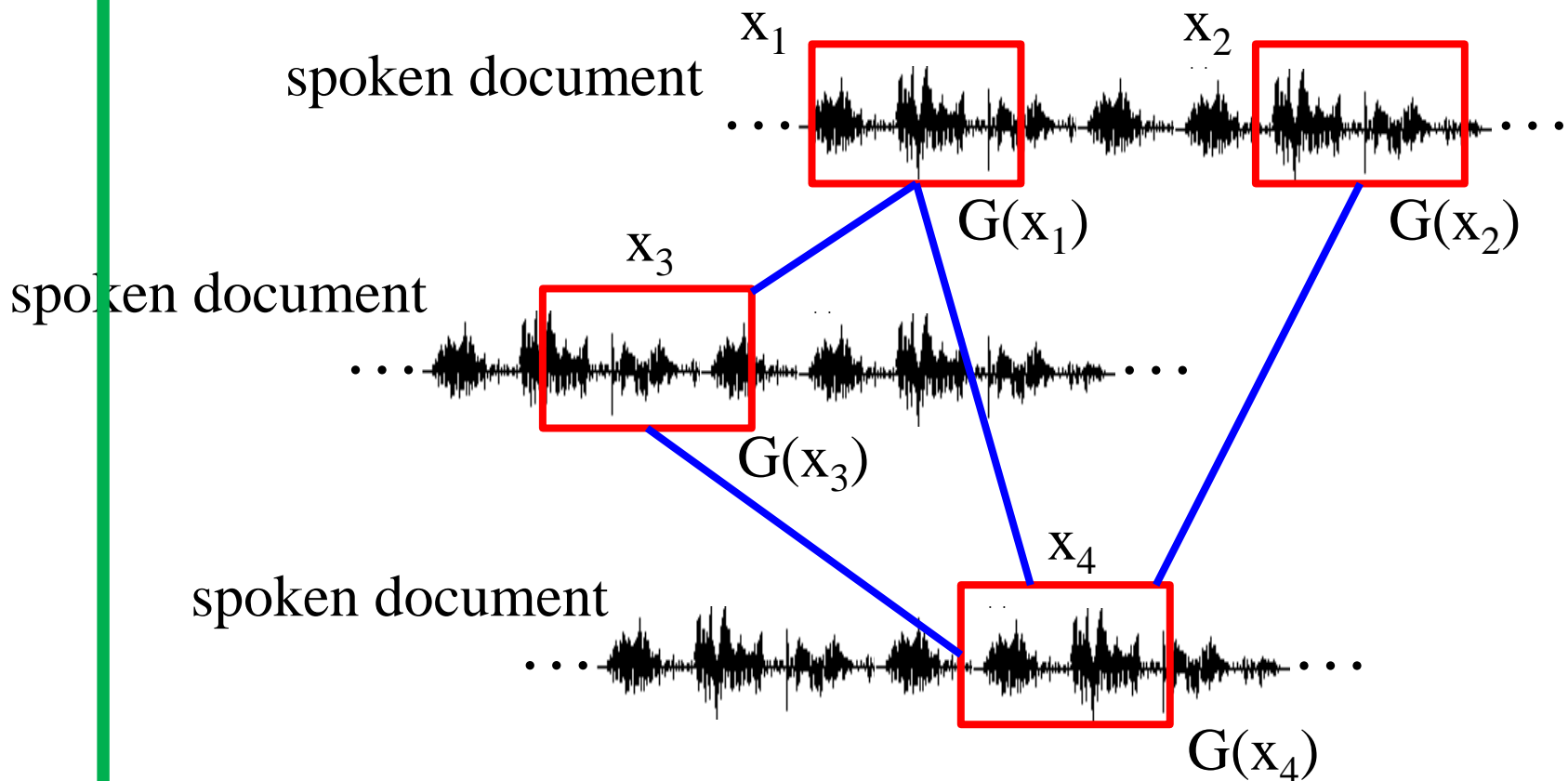
- For each word w in the lexicon



Obtain new score $G(x)$ by random walk

Graph-based Approach for Semantic Retrieval

- For each word w in the lexicon



Repeat this process for all the words w in the lexicon

Graph-based Approach for Semantic Retrieval

Lattice-based

Document Model:

$$P(w | \theta_d) = \frac{E(w, d)}{\sum_{w'} E(w', d)}$$

Graph-enhanced

document model:

$$P(w | \theta'_d) = \frac{E'(w, d)}{\sum_{w'} E'(w', d)}$$

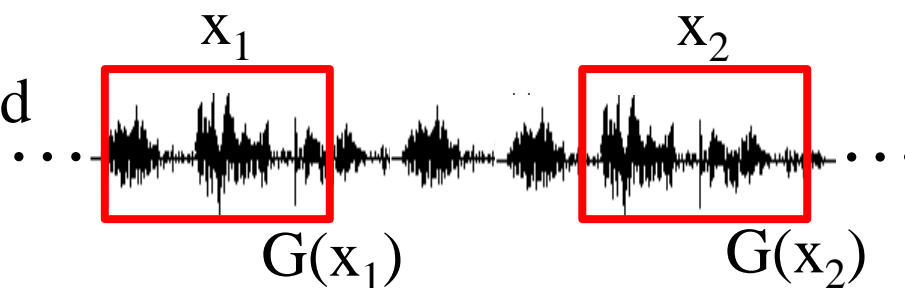
replace

Better estimation of term frequencies for each word w in d

$$E'(w, d)$$

query/document expansion can be applied

spoken document d

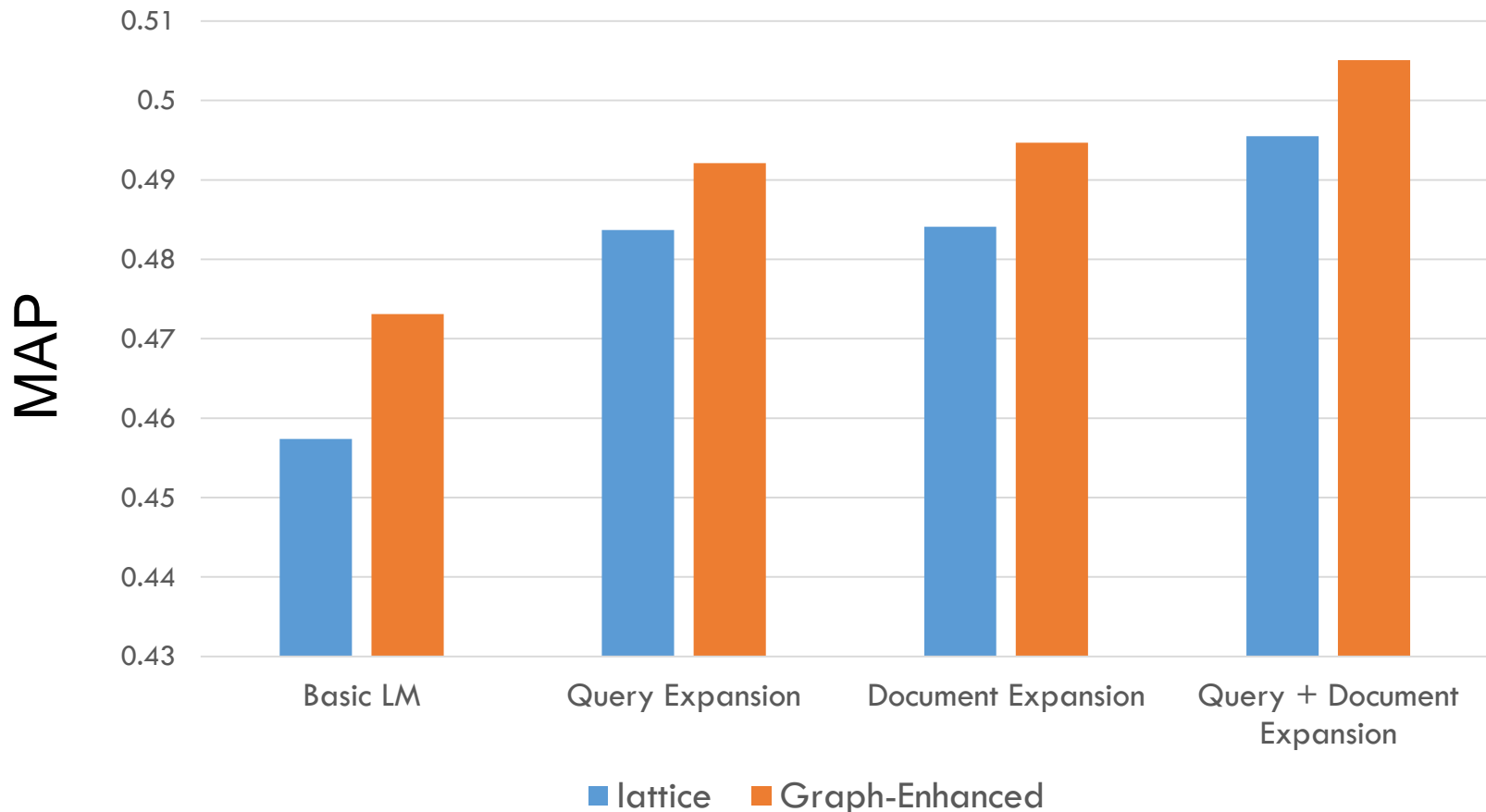


Scores from graph

Graph-based Approach for Semantic Retrieval - Experiments

□ Experiments on TV News

[Lee & Lee, IEEE/ACM T. ASL 14]

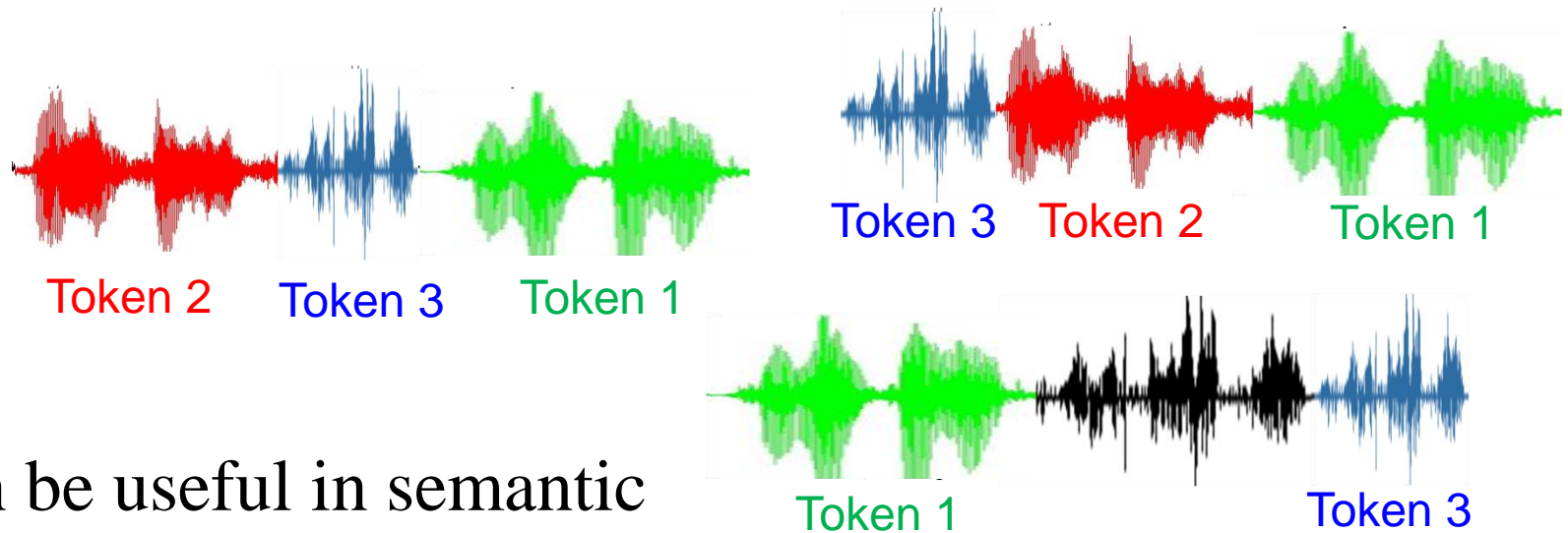


New Direction 4-2:
Special Semantic Retrieval Techniques
for Spoken Content
Exploiting Acoustic Tokens



Acoustic Tokens

- We can discover “*acoustic tokens*” in direction 3



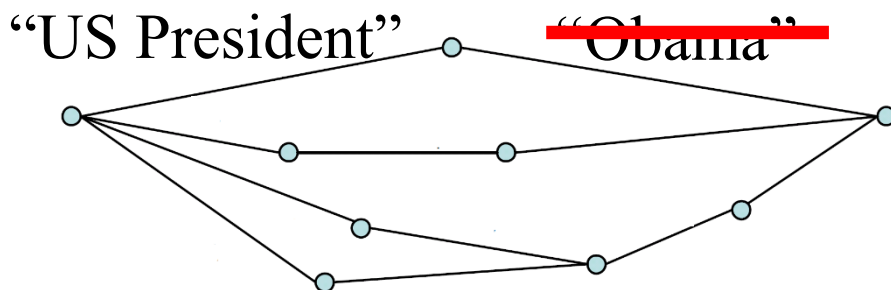
Can be useful in semantic retrieval of spoken content:

- Query expansion with acoustic tokens
- Unsupervised semantic retrieval of spoken content

Query Expansion with Acoustic Tokens

Basic idea of query expansion

Related terms frequently co-occur
in the same spoken document



If “Obama” is not in the lexicon

➔ “Obama” will never appear in lattices.

➔ We can never know “Obama” co-occur with
“US President” in query expansion.

➔ Typical approach: using subwords

➔ Query expansion with acoustic
tokens

} Complementary
to each other

Query Expansion with Acoustic Tokens

[Lee & Lee, ICASSP 13]

Original Text Query:

“US President”



d_{100} : US President ...

d_{205} : ... US President

First pass: Retrieve spoken documents containing “US President” in the transcriptions

Query Expansion with Acoustic Tokens

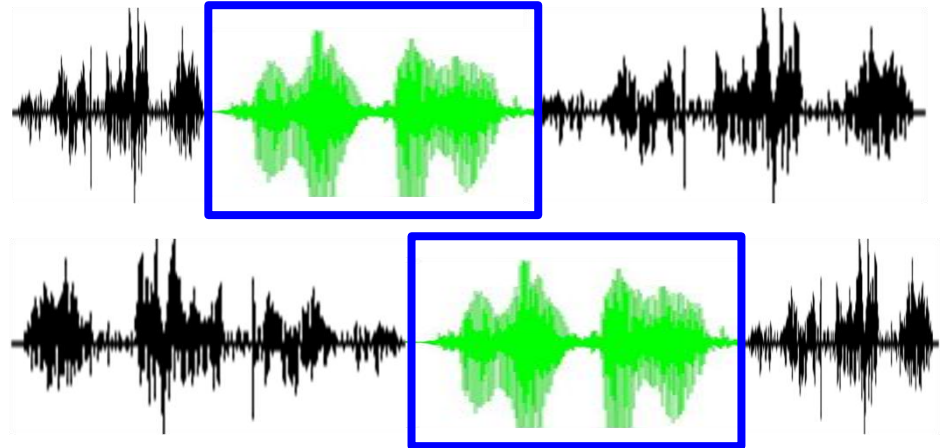
Original Text Query:

“US President”



d_{100} : US President ...

d_{205} : ... US President



Find acoustic tokens frequently appear in the signals of these retrieved documents

Query Expansion with Acoustic Tokens

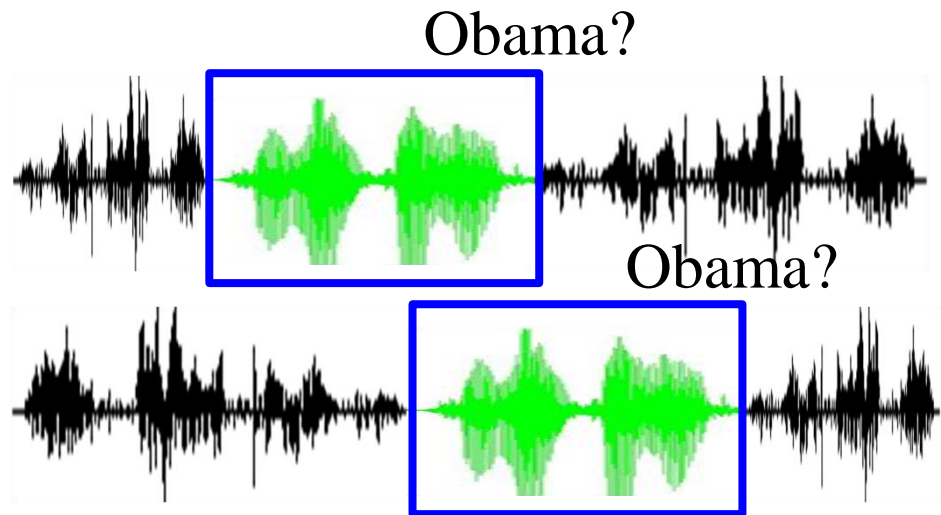
Original Text Query:

“US President”



d_{100} : US President ...

d_{205} : ... US President



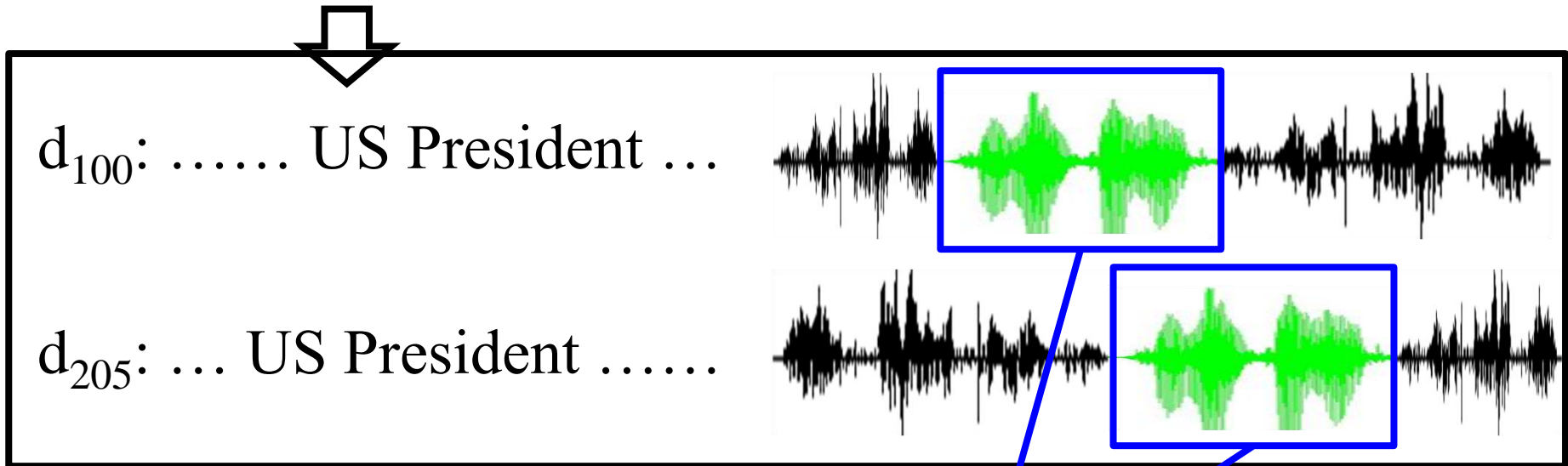
Even the terms related to the query is OOV

- ➡ If they co-occur with the query in speech signals
- ➡ Find acoustic tokens corresponding to these terms

Query Expansion with Acoustic Tokens

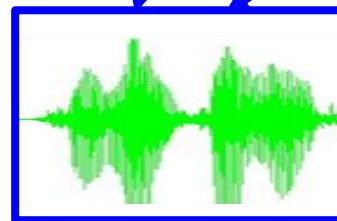
Original Text Query:

“US President”

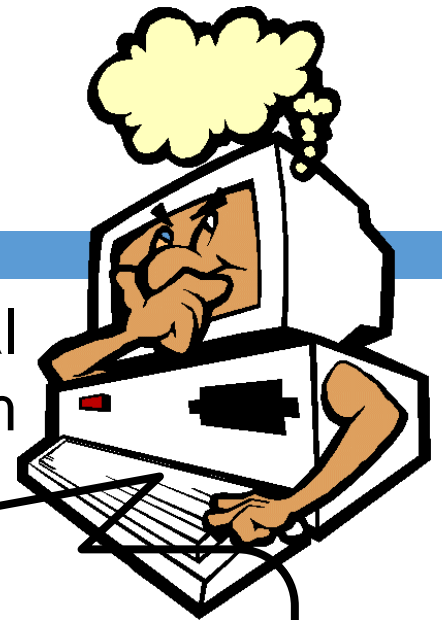


Expanded Query:

“US President” +



Query Expansion with Acoustic Tokens



Retrieval system



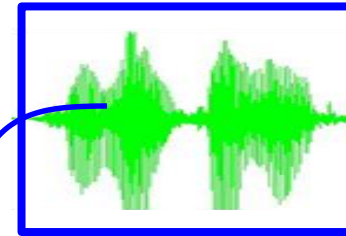
user

“US President”

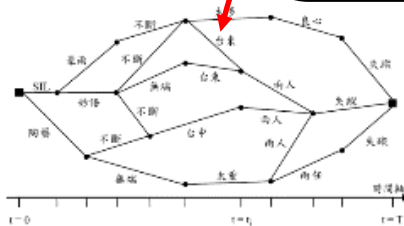
Expanded Query:

“US President”
“White House”

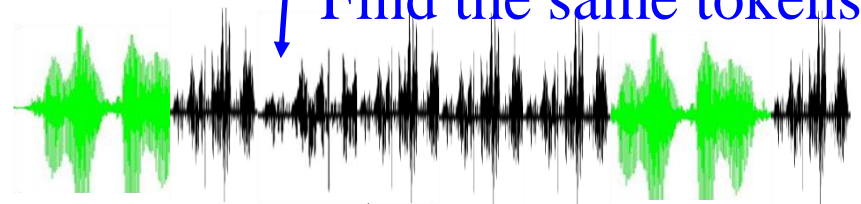
and



Find the same tokens



Lattices



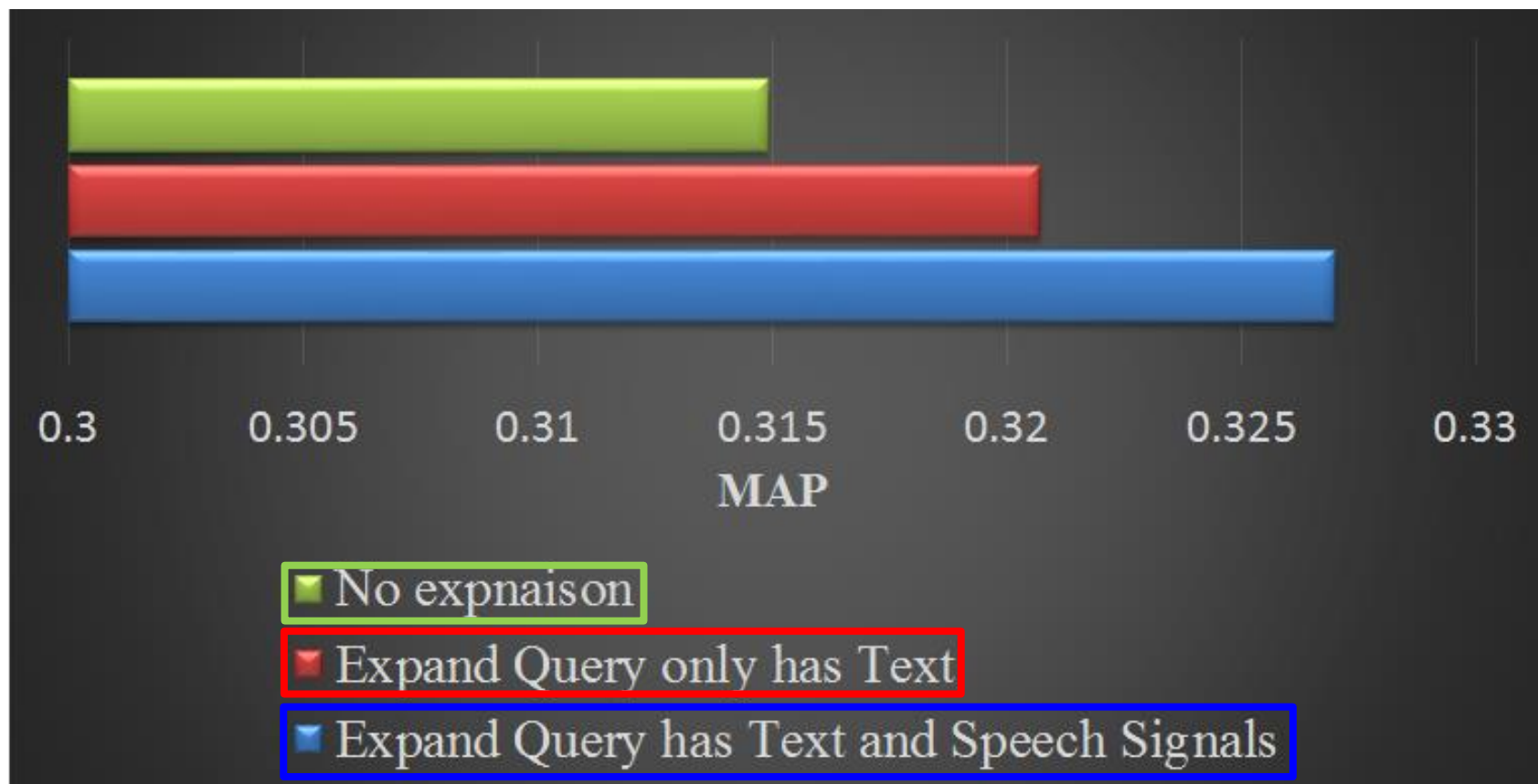
By expanding the text query with acoustic tokens, more semantically related audio files can be retrieved.

Query Expansion

– Acoustic Patterns

□ Experiments on TV News

[Lee & Lee, ICASSP 13]



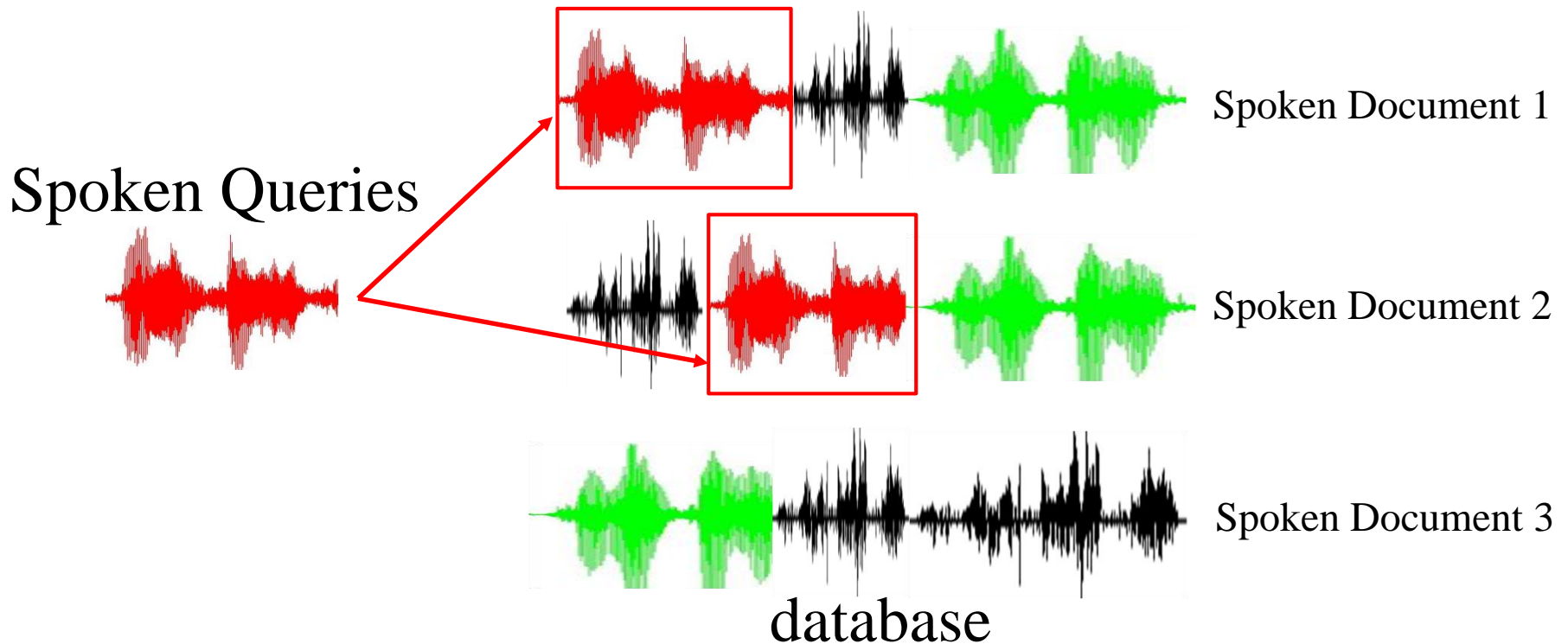
Unsupervised Semantic Retrieval

- *Unsupervised Semantic Retrieval* [Li & Lee, ASRU 13][Oard, FIRE 13]
 - ▣ Find spoken documents *semantically related* to the *spoken queries*
 - ▣ **Without speech recognition**
- New task, not too much previous work
 - ▣ Below is just a very preliminary study based on query expansion with acoustic tokens [Li & Lee, ASRU 13]

Unsupervised Semantic Retrieval

1. Find spoken documents containing the spoken query

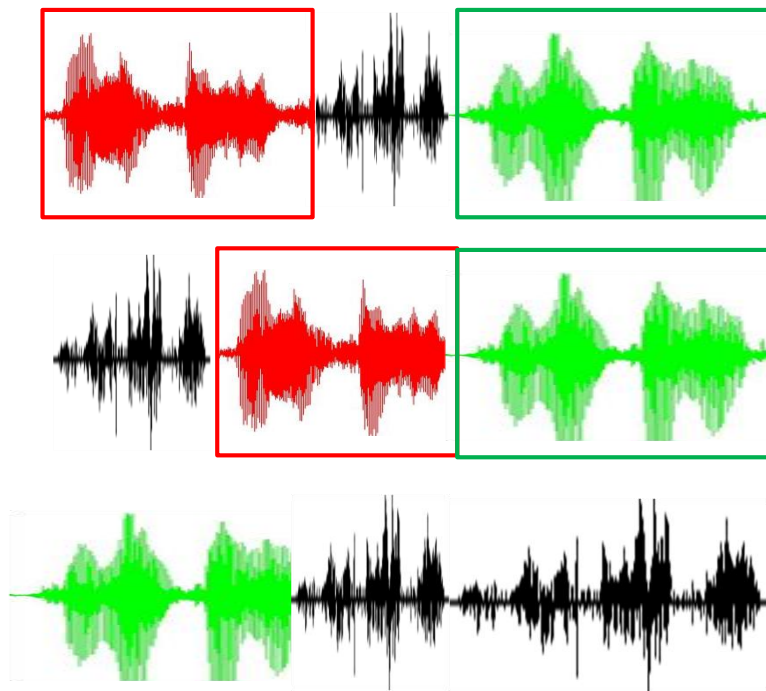
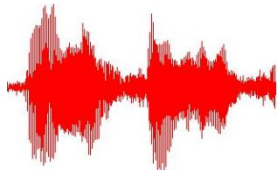
➔ Done by the query-by-example spoken term detection approaches (e.g. DTW)



Unsupervised Semantic Retrieval

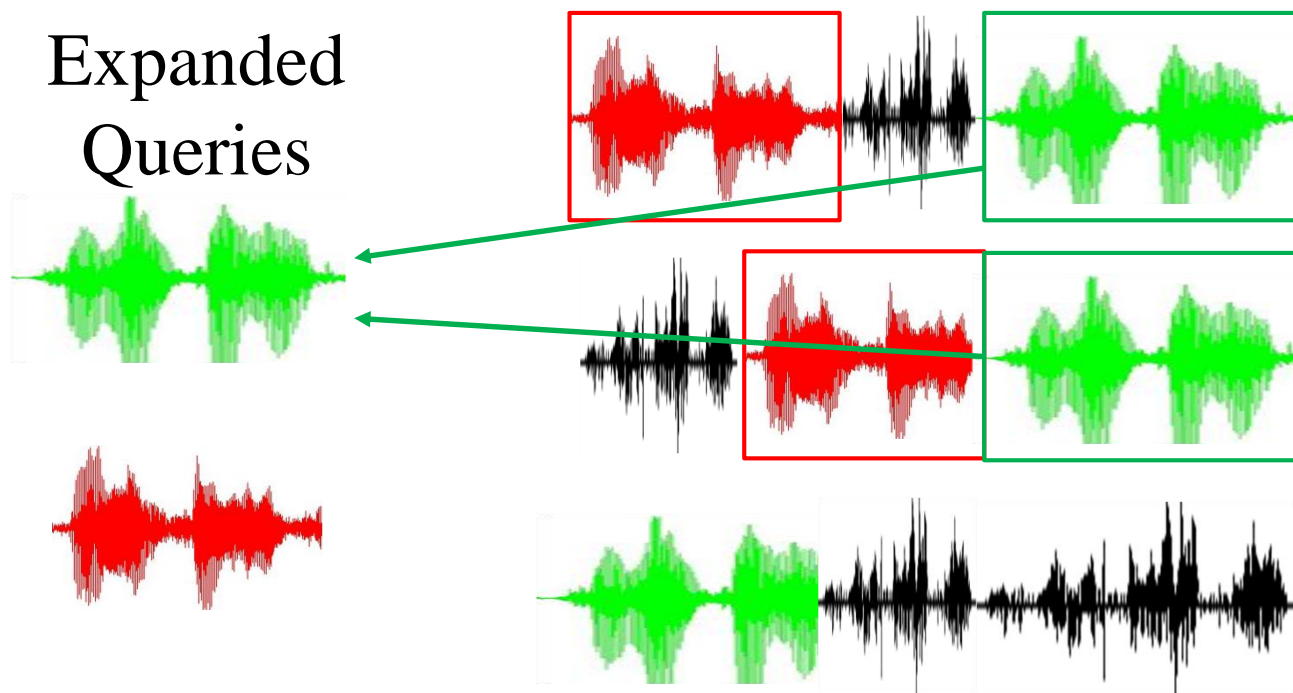
2. Find acoustic tokens frequently co-occurring with the spoken queries in the same document

Spoken Queries



Unsupervised Semantic Retrieval

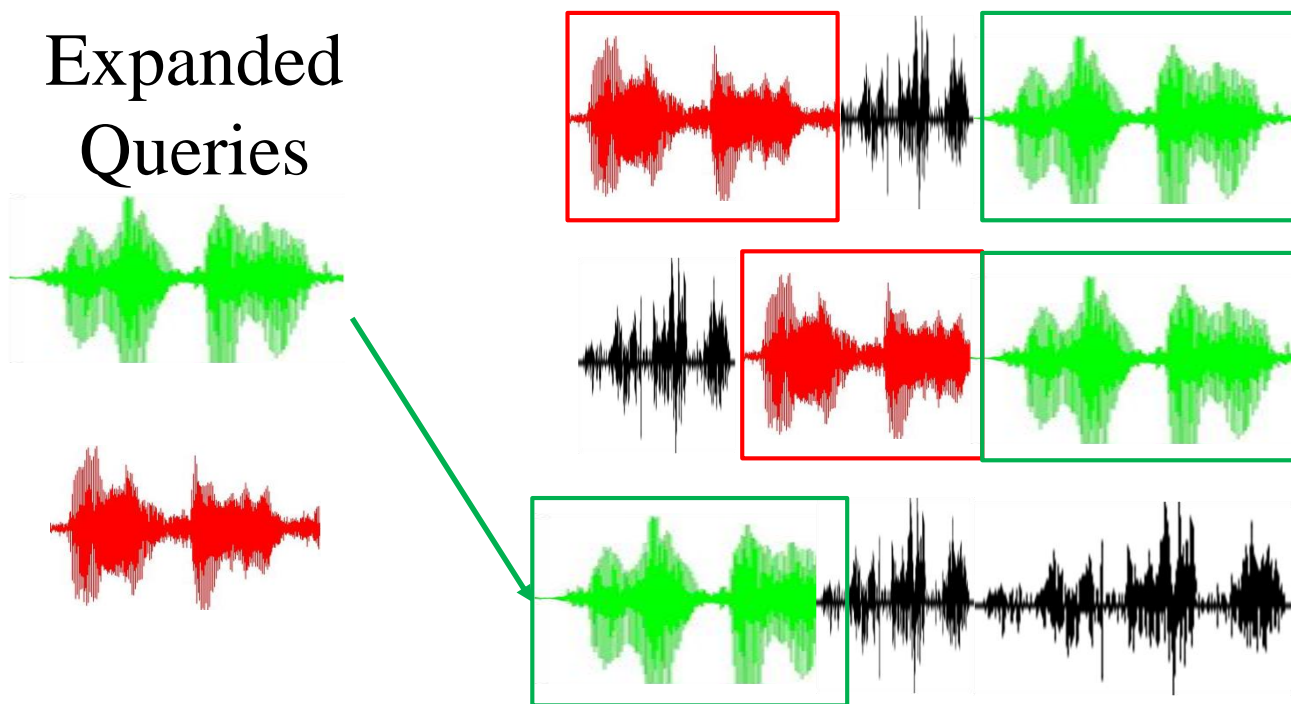
3. Use the acoustic tokens to expand the original spoken query



Unsupervised Semantic Retrieval

4. Retrieve again by the expanded queries

➔ Can retrieve spoken documents not containing the original spoken queries



Unsupervised Semantic Retrieval

- Experiments

- Broadcast news, MAP as evaluation measure
 - ▣ Using only DTW for unsupervised semantic retrieval: MAP= 8.76%
 - The semantically related documents without the query term cannot be retrieved by DTW.
 - ▣ Expanded by Acoustic Tokens: MAP= 9.70%
 - Unsupervised semantic retrieval has a long way to go

New Direction 5: Speech Content is Difficult to Browse!

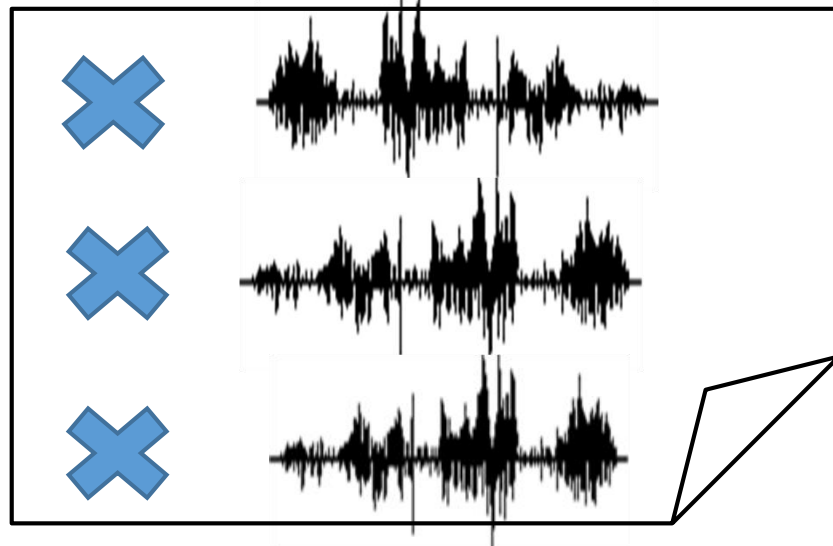


Audio is hard to browse

- When the system returns the retrieval results, user doesn't know what he/she get at the first glance



Retrieval Result



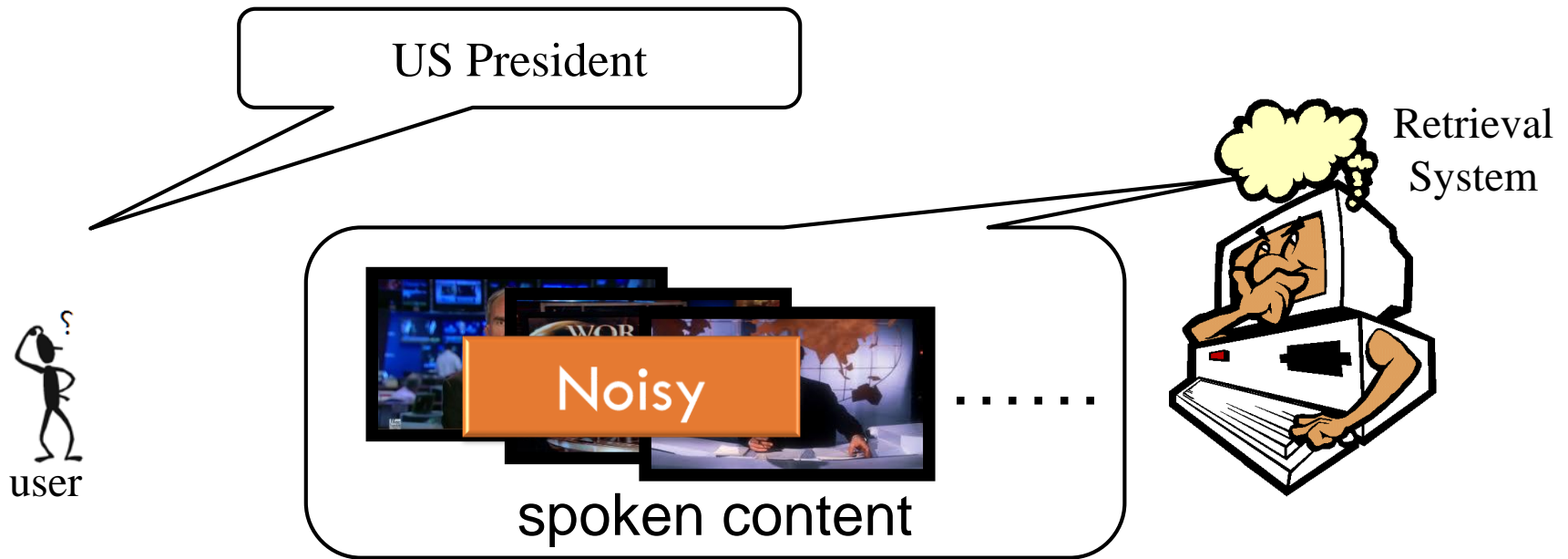
Audio is hard to browse

- ❑ Interactive spoken content retrieval
- ❑ Extracting Core Information
- ❑ Organizing Retrieved Results
- ❑ Spoken Question answering

New Direction 5-1:
Speech Content is Difficult to Browse!
Interactive Spoken Content Retrieval



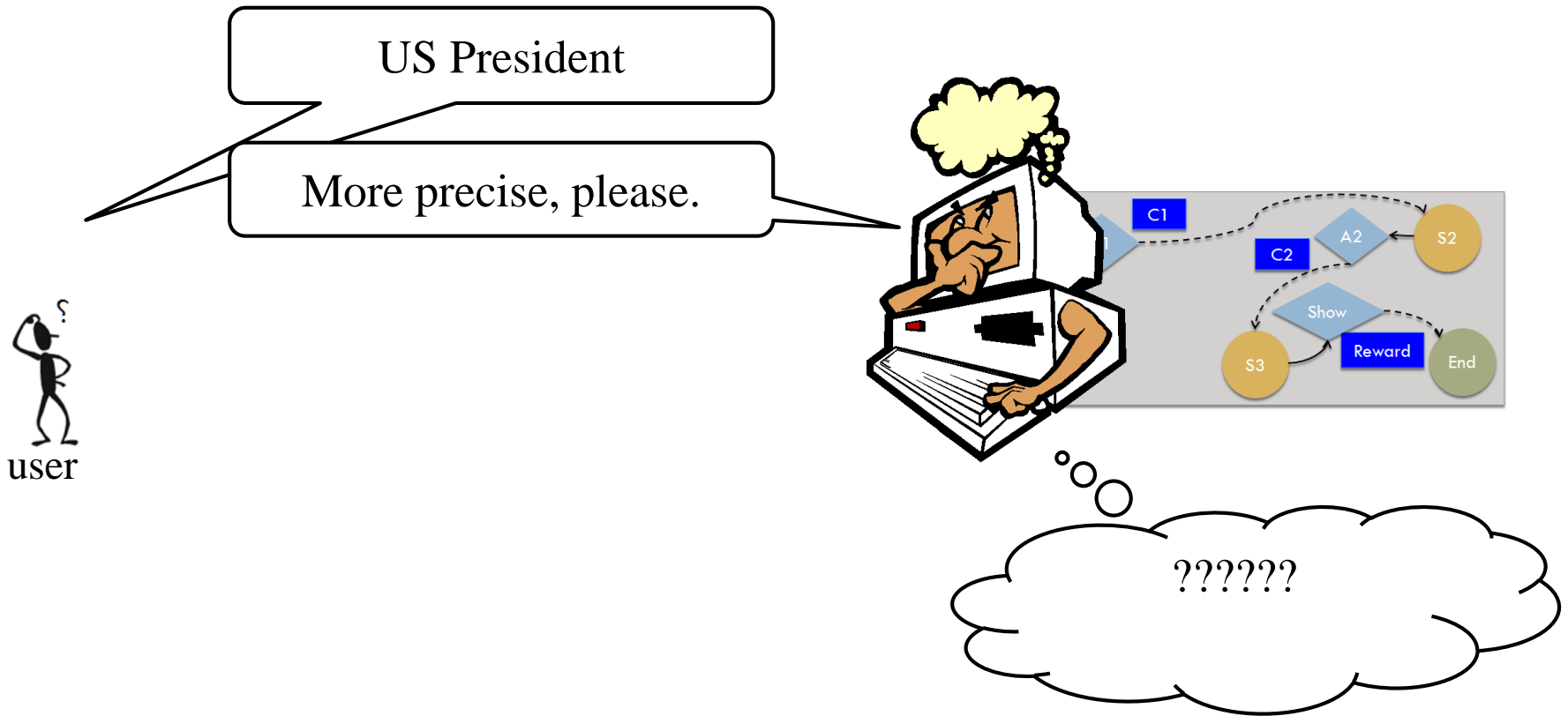
Interactive spoken content retrieval



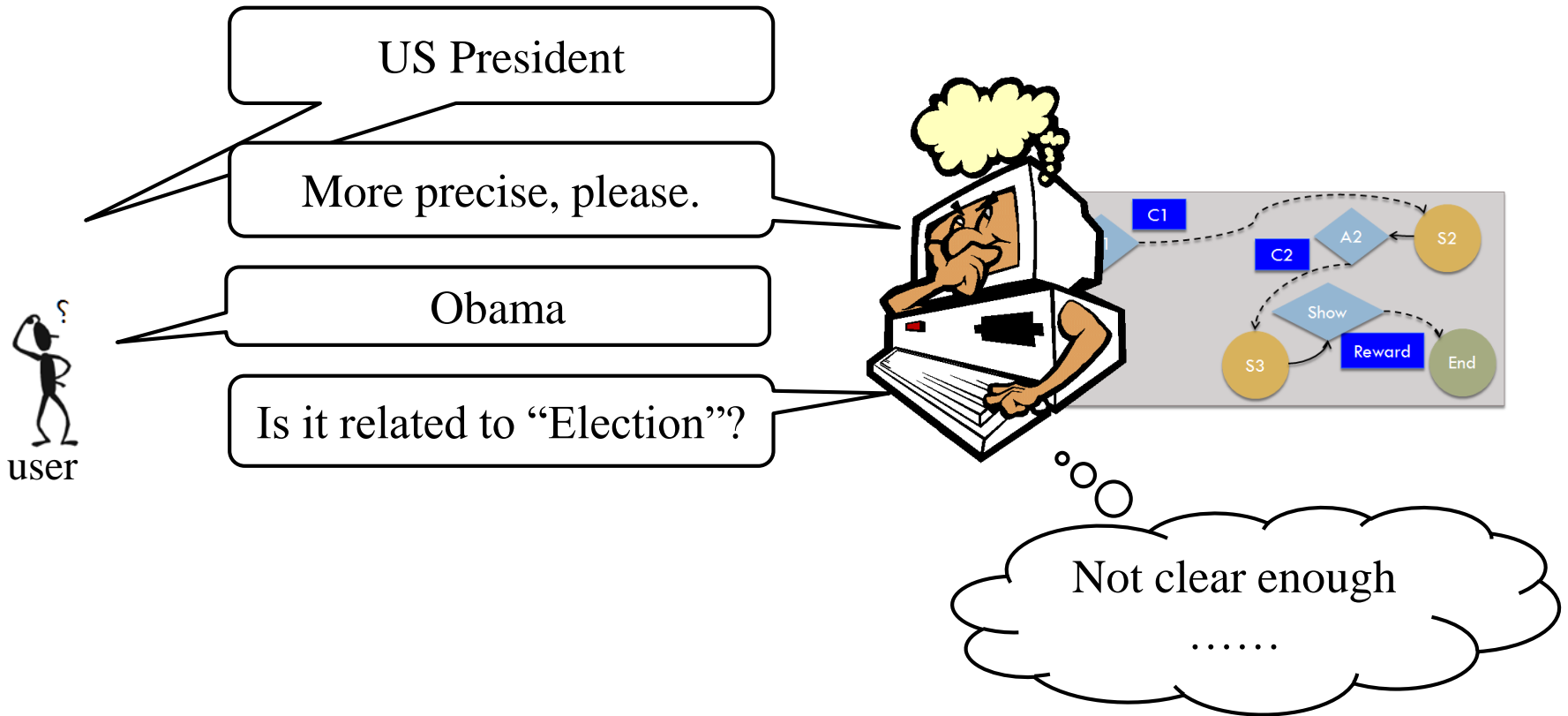
Input query is usually short → Cannot describe the information need clearly

Speech recognition always produces errors.

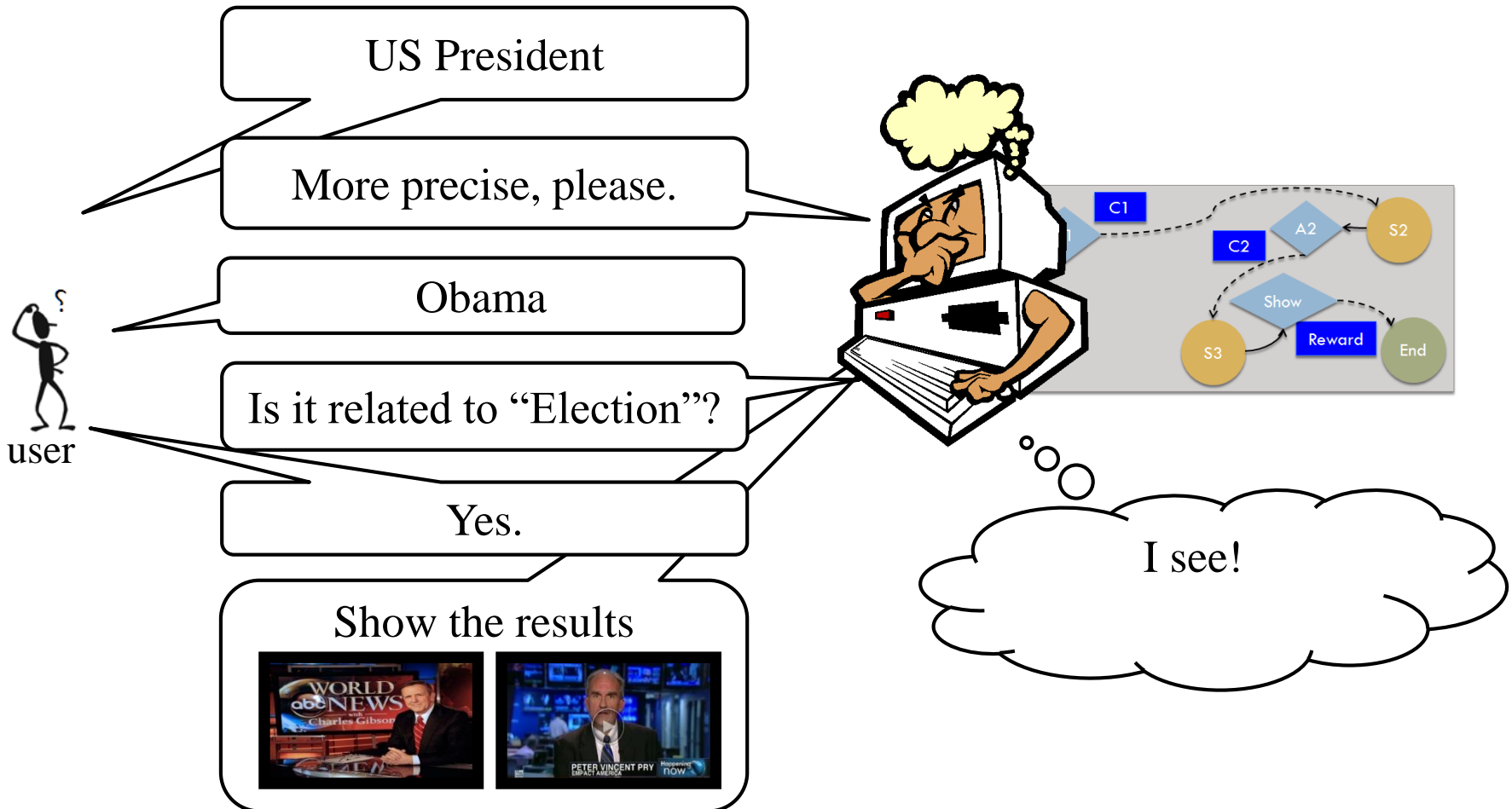
Interactive spoken content retrieval



Interactive spoken content retrieval

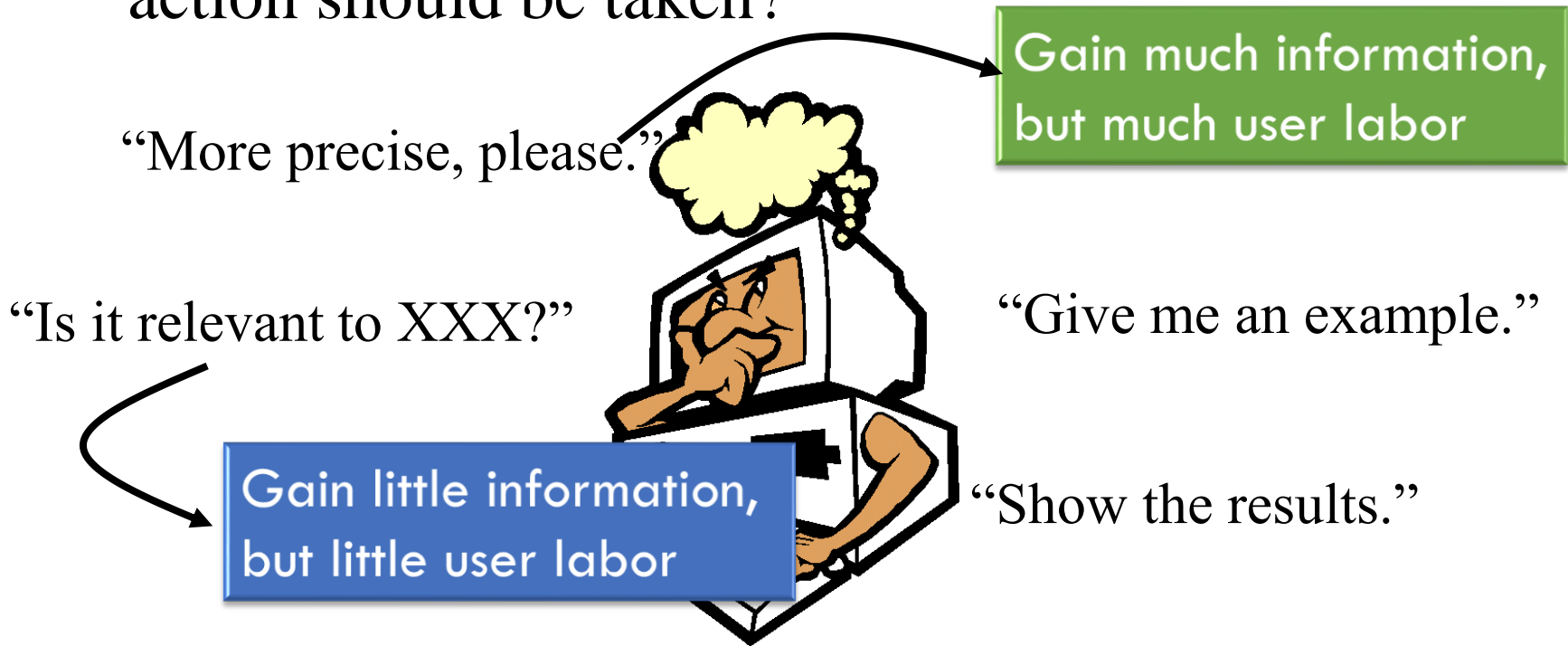


Interactive spoken content retrieval



Challenges

- Given the information entered by the users, which action should be taken?



Borrowing the experiences from developing dialogue system (air ticket booking, city guides, personal assistant ..)

MDP for Interactive Retrieval

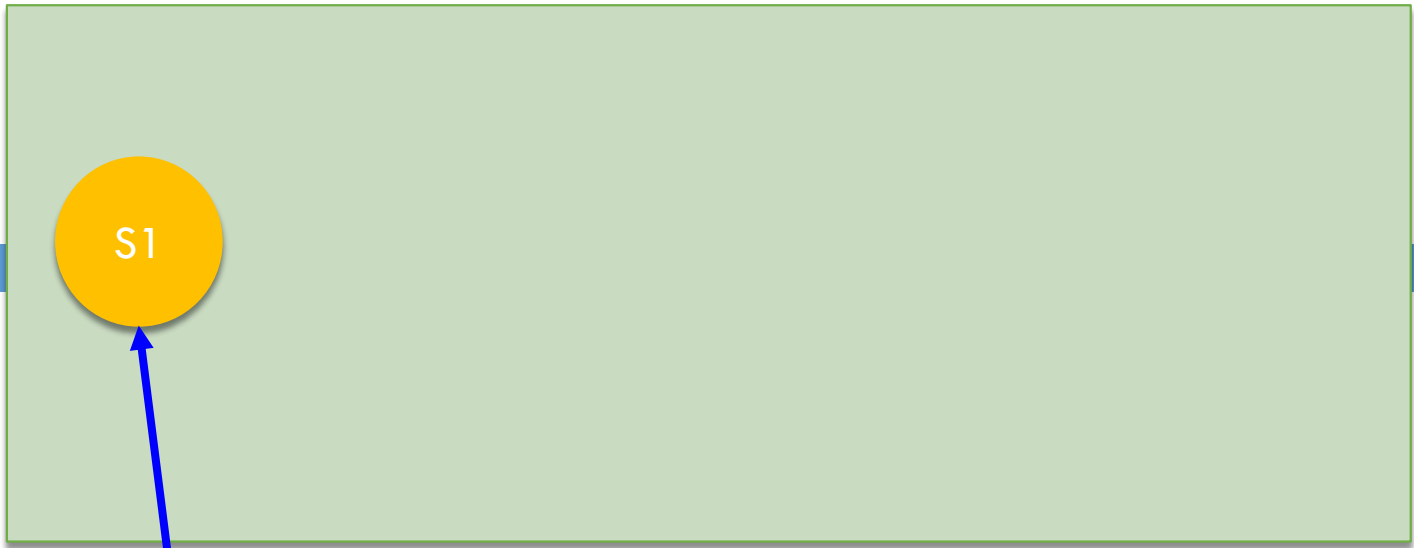
- Markov Decision Process (MDP)
 - ▣ The system is in certain states.
 - ▣ Which action should be taken depends on the state the system is in.
- MDP for Interactive retrieval [Wen & Lee, Interspeech 12][Wen & Lee, ICASSP 13]
 - ▣ *State*: the degree of clarity of the user's information need

Ambiguous

state space

Clear

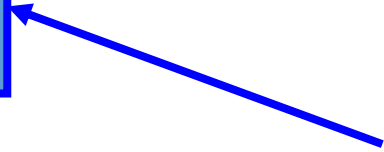
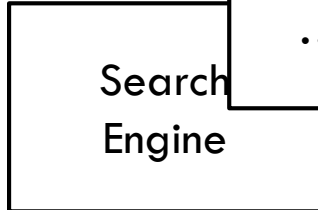
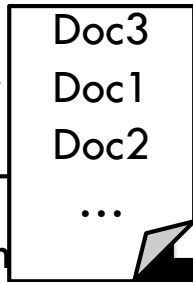


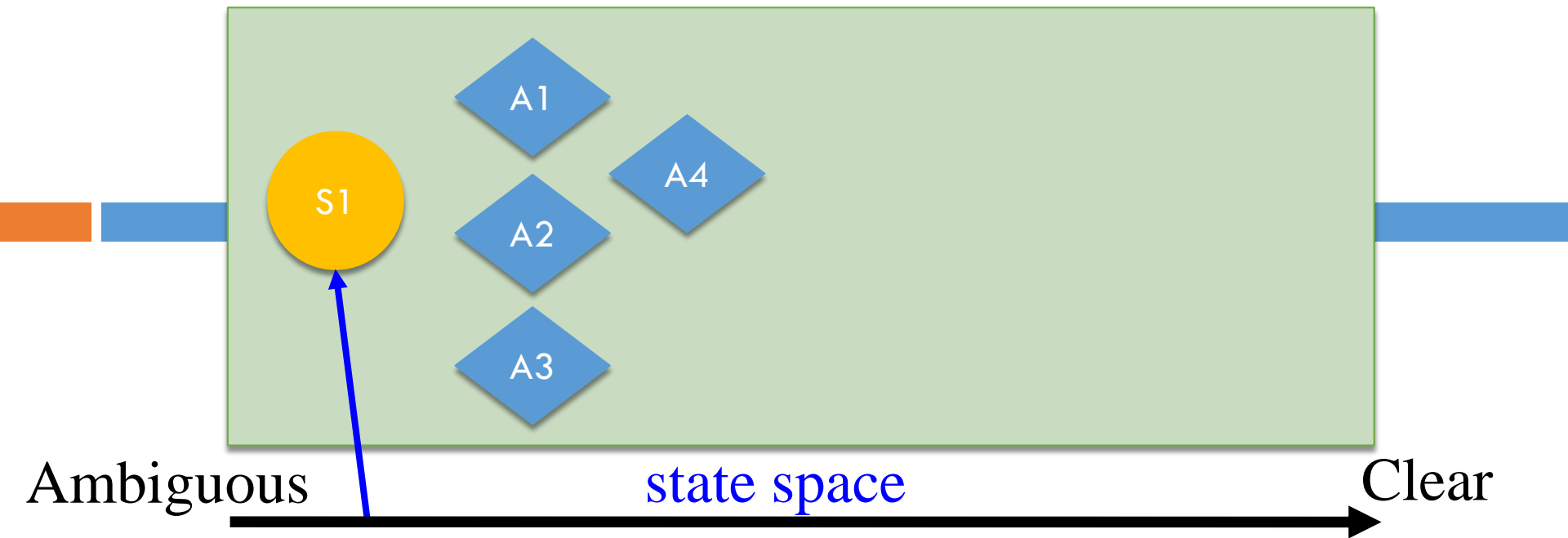


Ambiguous state space Clear

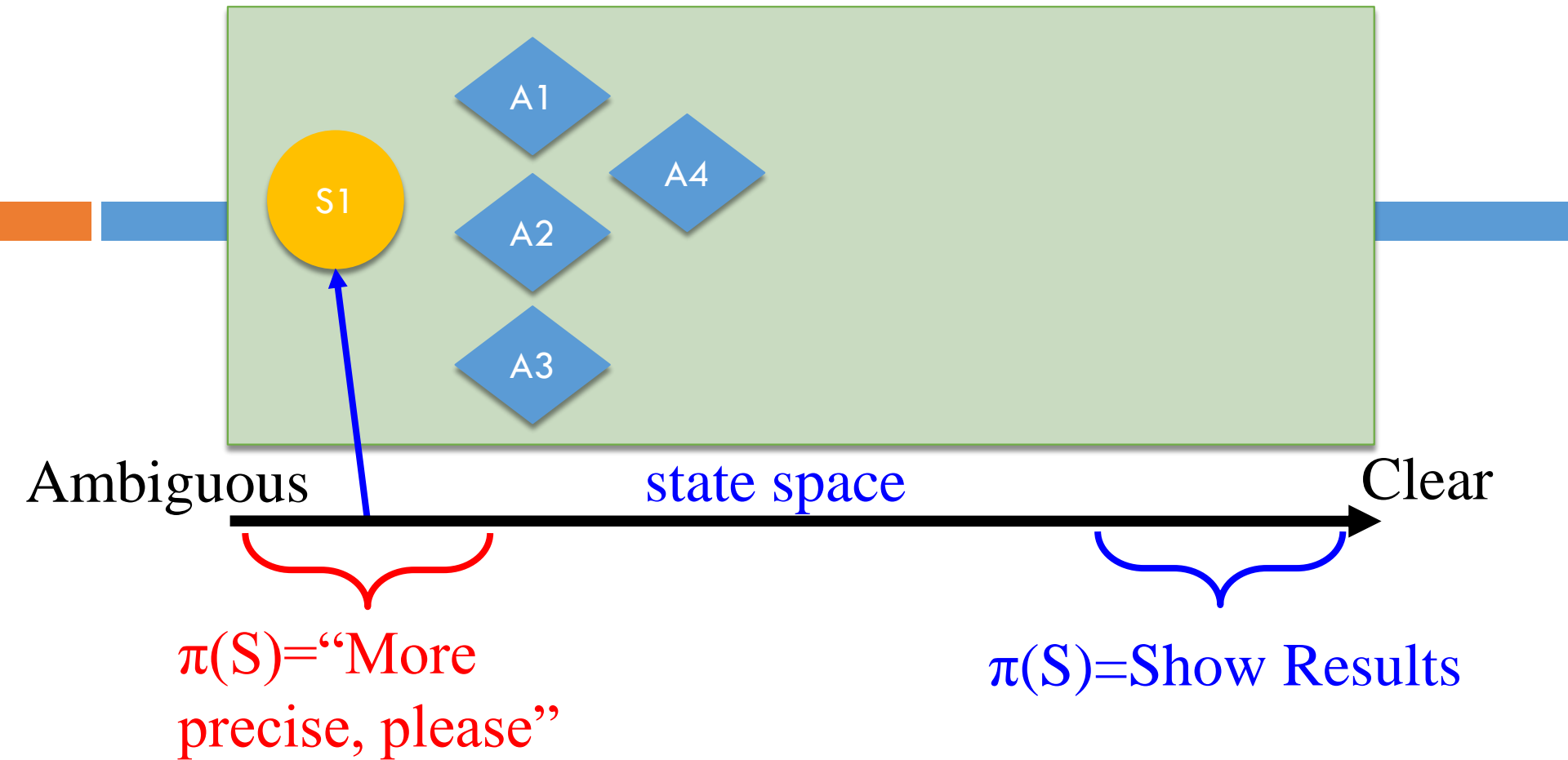
State Estimator: Estimate the degree of clarity from the retrieval results

[Cronen-Townsen, SIGIR 02]
[Zhou, SIGIR 07]

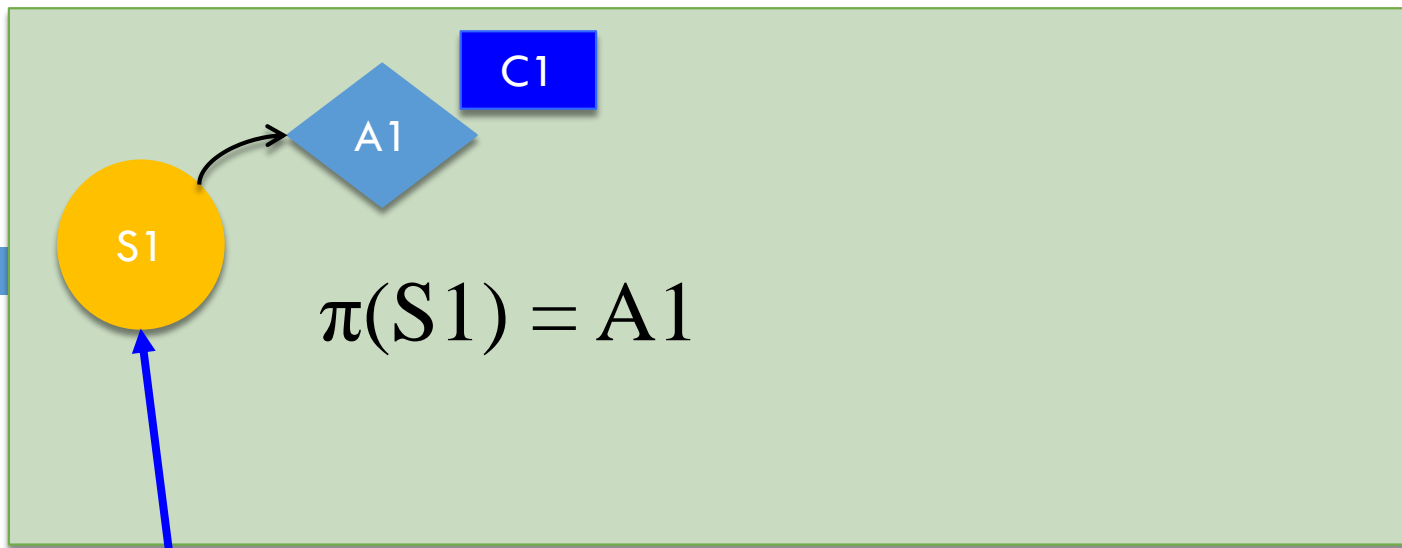




- A set of candidate actions
 - ▣ System: “More precisely, please.”
 - ▣ System: “Is it relevant to XXX?”
 - ▣
- There is an action “show results”
 - ▣ When the system decides to show the results, the retrieval session is ended



- Choose the actions by intrinsic policy $\pi(S)$
 - The policy is a function
 - Input: state S , output: action A



Ambiguous

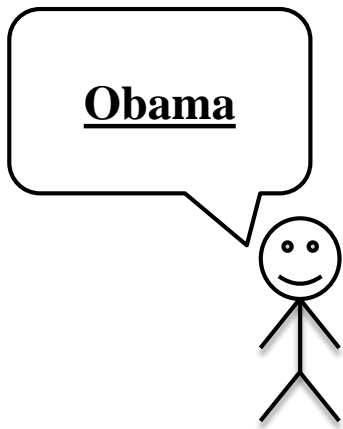
state space

Clear

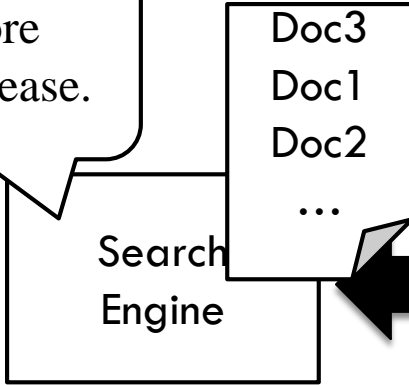
User response

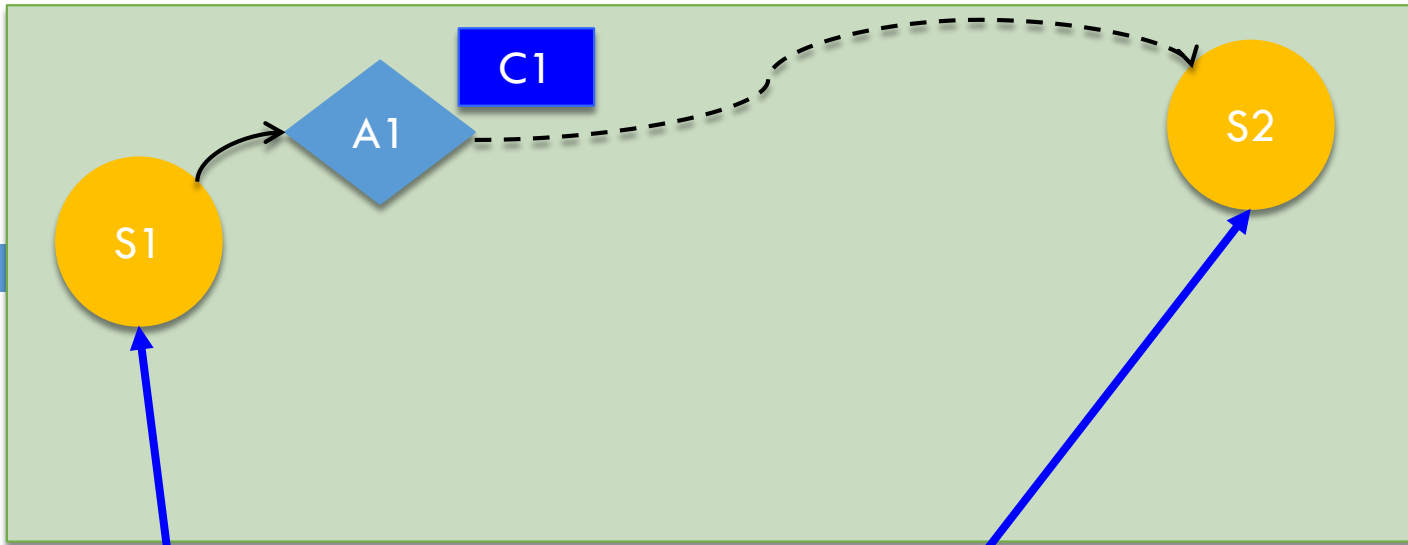


The system gets a cost C1 due to user labor.



A1: More precise, please.



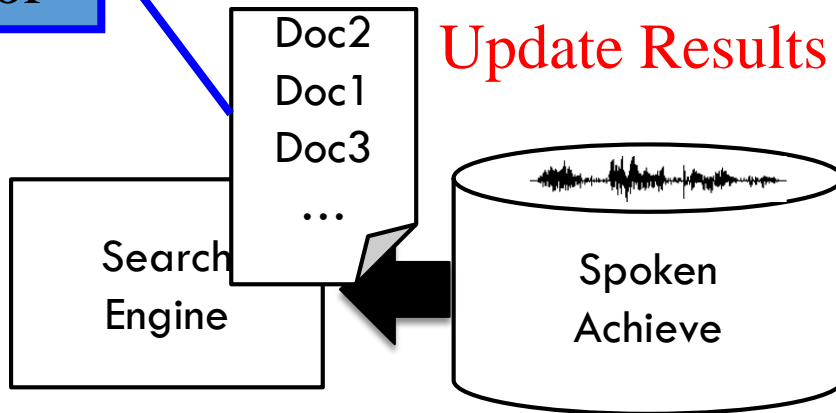


Ambiguous

state space

Clear

State Estimator



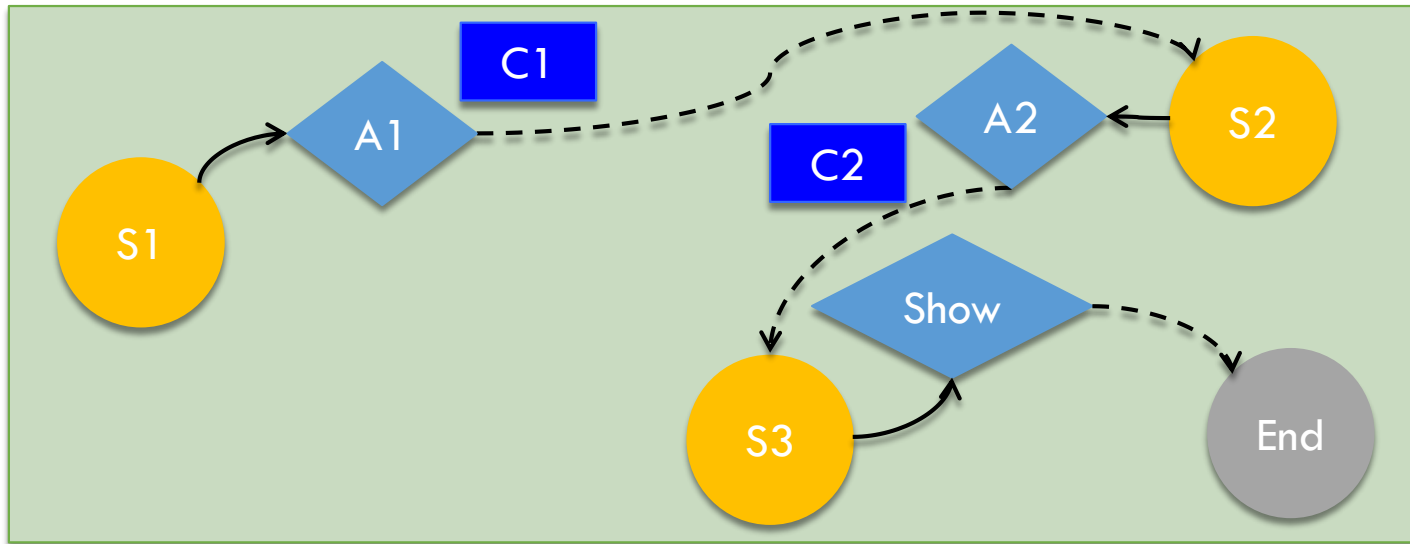
Update Results

Search Engine

Doc2
Doc1
Doc3
...

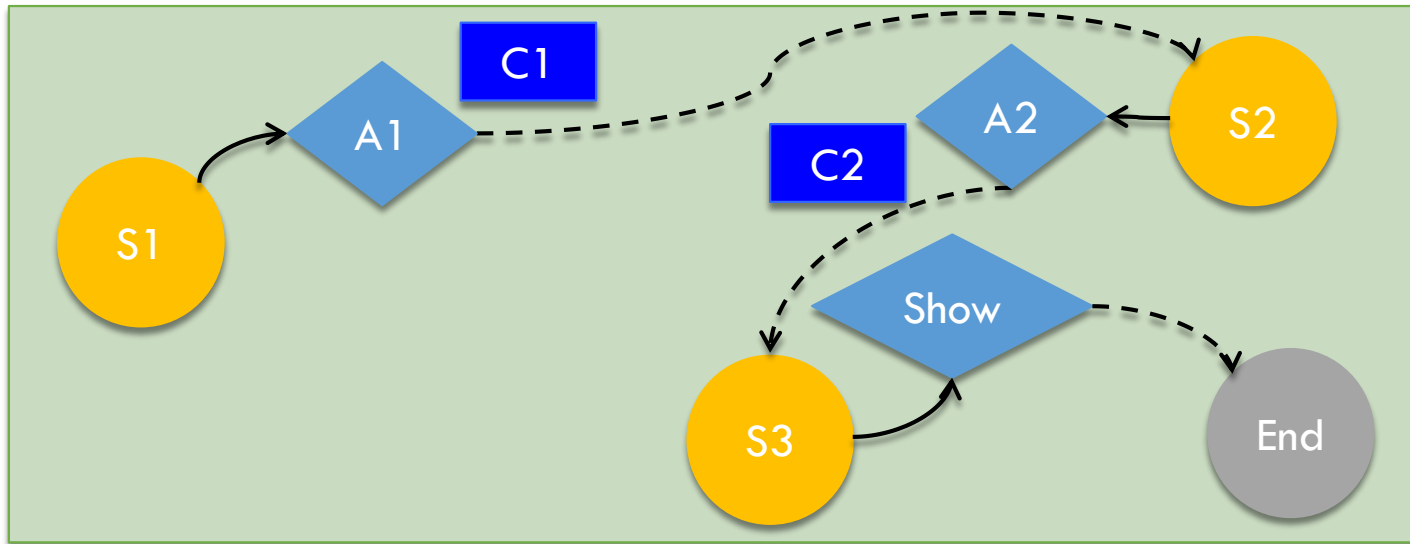
Spoken Achieve

Interact with Users - MDP



- Good interaction:
 - ▣ The quality of final retrieval results shown to the users are as good as possible
 - ▣ The user labors (C1, C2) are as small as possible

Interact with Users - MDP



- Find a policy π that

Maximizing Retrieval Quality, Minimizing User Labor

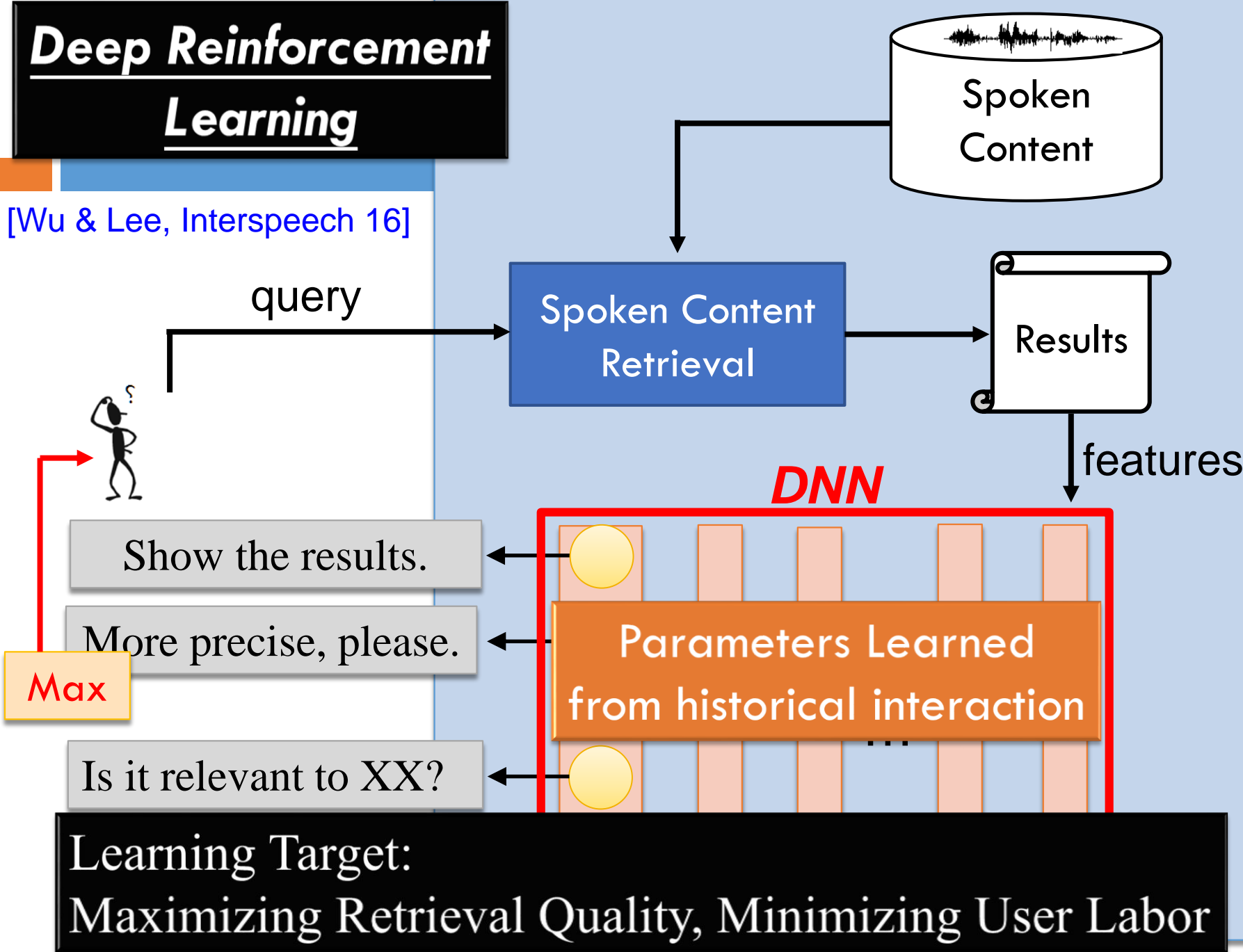
- The policy π can be learned from historical interaction [Wen & Lee, Interspeech 12][Wen & Lee, ICASSP 13]

Deep Reinforcement Learning



Deep Reinforcement Learning

[Wu & Lee, Interspeech 16]



Max

Show the results.
More precise, please.
Is it relevant to XX?

Learning Target:
Maximizing Retrieval Quality, Minimizing User Labor

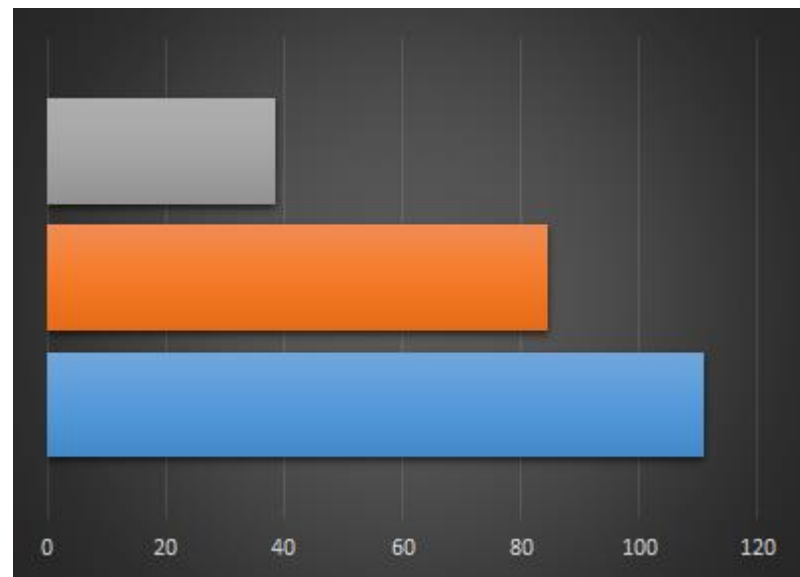
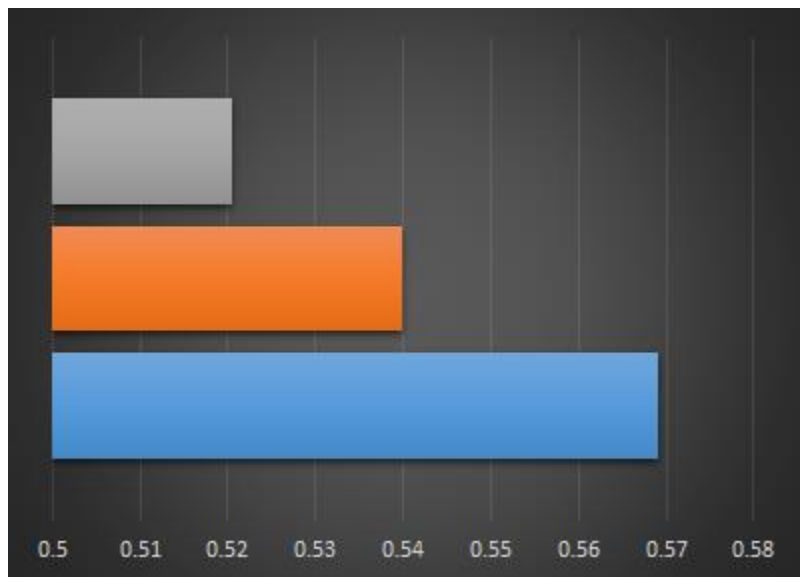
Experimental Results

- Broadcast news, semantic retrieval

Optimization Target:

Retrieval Quality (MAP)

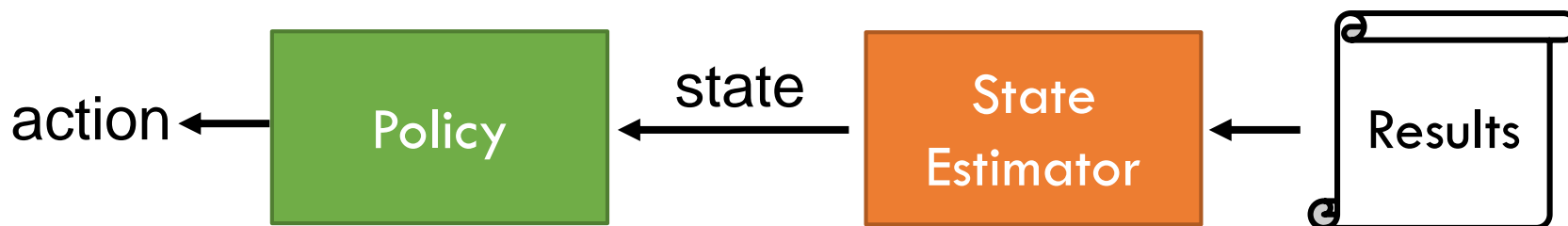
Retrieval Quality - User labor



■ Hand-crafted ■ MDP ■ Deep Reinforcement Learning

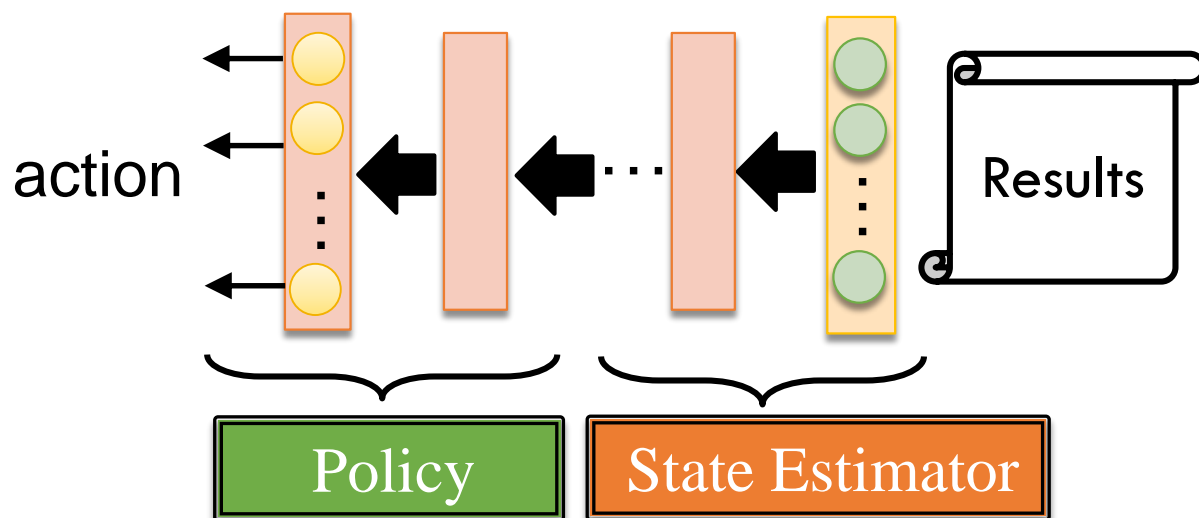
Deep Reinforcement Learning v.s. MDP for interactive retrieval

- MDP for interactive retrieval [Wen & Lee, Interspeech 12][Wen & Lee, ICASSP 13]



The two stages were learned separately.

- Deep
End-to-end learning

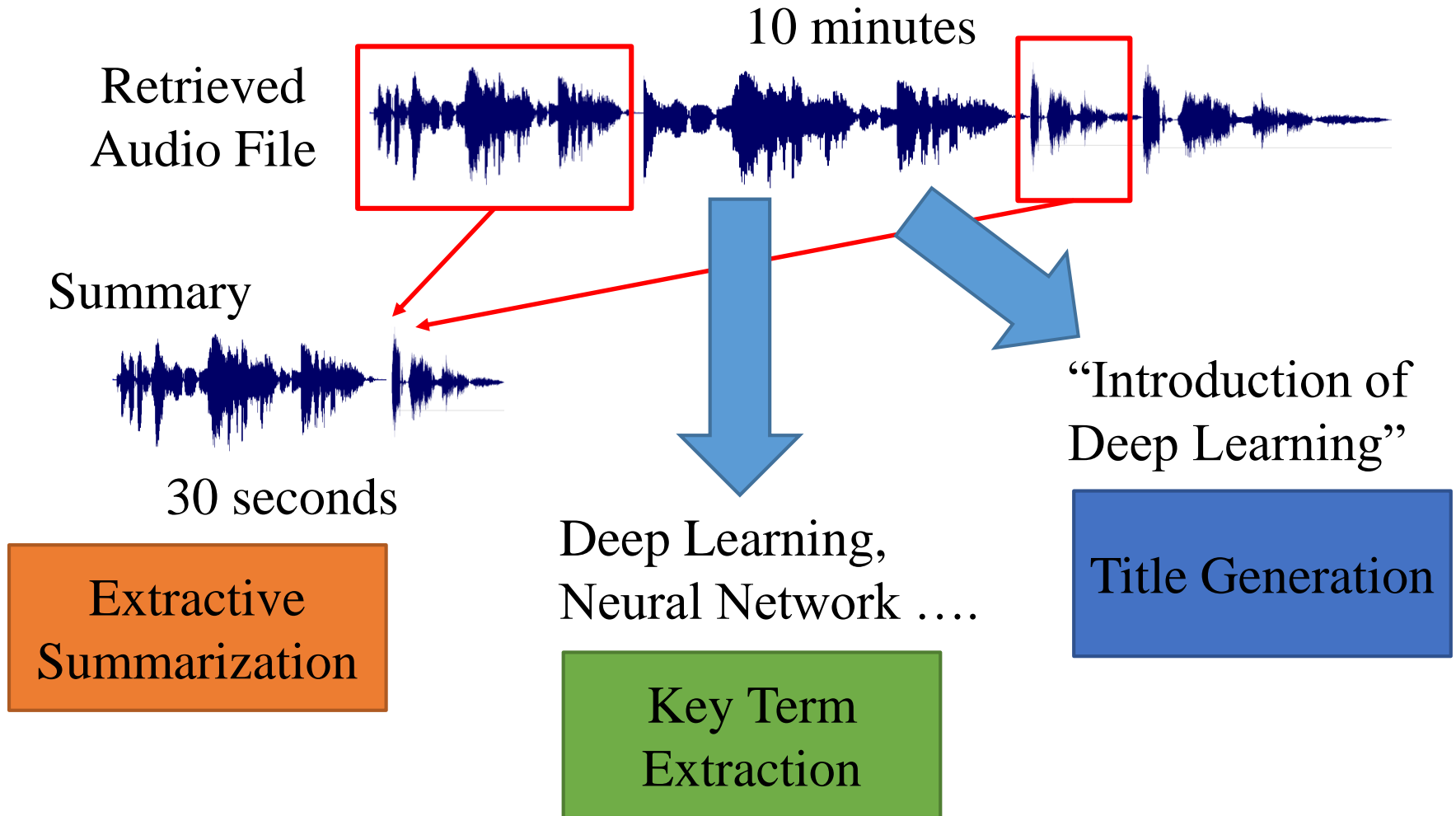


No hand-crafted states

New Direction 5-2:
Speech Content is Difficult to Browse!
Extracting Core Information



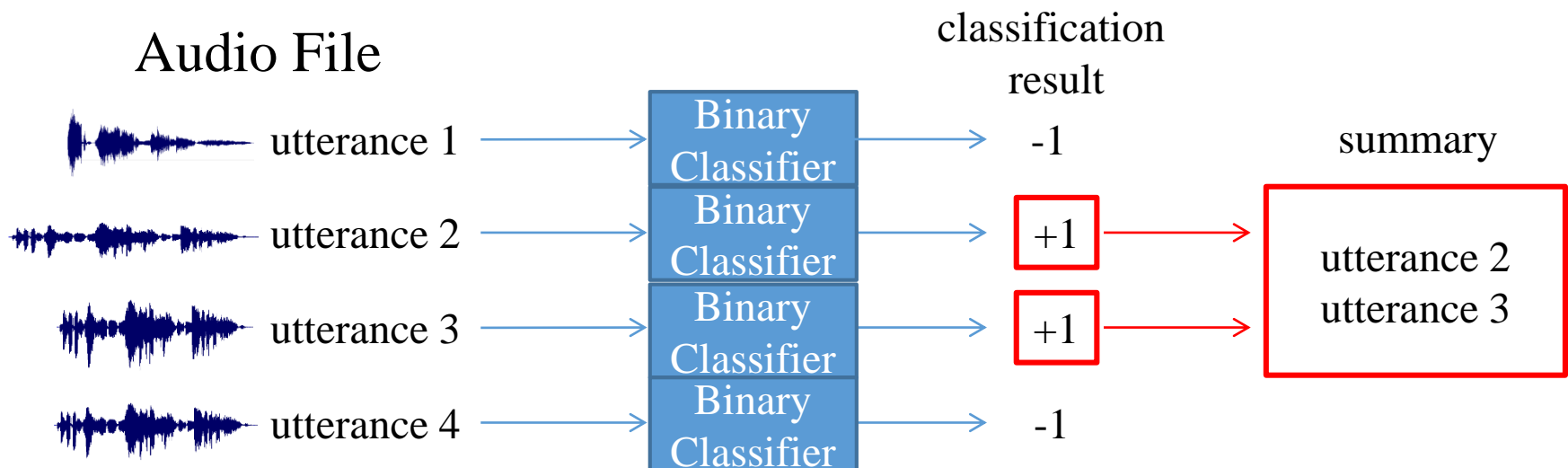
Extracting Core Information



Summarization

Reference: **13 Speech Summarization**
(Gokhan Tur, Renato De Mori, Yang Liu, Dilek Hakkani-Tür). G. Tur and R. DeMori, Spoken Language Understanding: Systems for Extracting Semantic Information from Speech.

- **Unsupervised Approach: Maximum Margin Relevance (MMR) and Graph-based Approach**
- **Supervised approach**
 - Naïve approach: Summarization problem can be formulated as binary classification



Summarization

– Binary Classification

- Binary classifier individually considers each utterance
- Not sufficient
 - ▣ Example: summary should be concise

Lecture Recording

Hello

LSA is Latent semantic analysis

LSA is useful for summarization

Therefore

LSA improves summarization

.....

Summary

LSA is useful for summarization
LSA improves summarization

Summary should be succinct

To generate a good summary, “*global information*” should be considered

More advanced machine learning techniques

Summarization

- Considering Global Information

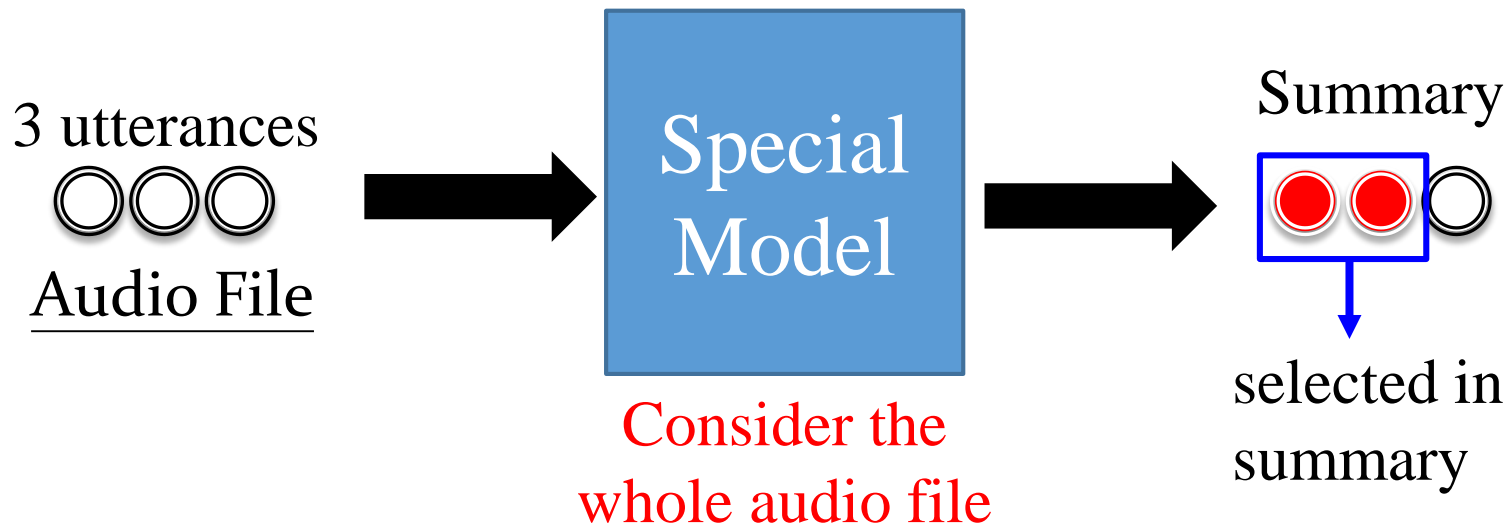
- Learn a special model

- ▣ Input: whole audio file

- ▣ Output: summary

[Lee & Lee, ICASSP 13]

[Lee & Lee, Interspeech 12]



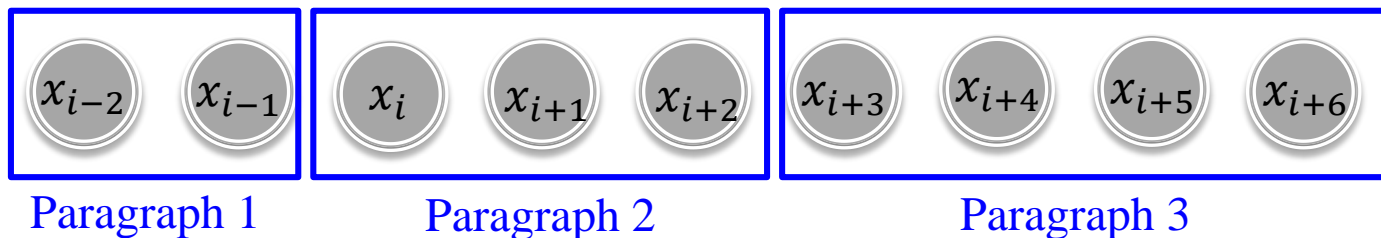
Structured SVM: I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun. Support Vector Learning for Interdependent and Structured Output Spaces, ICML, 2004.

Summarization

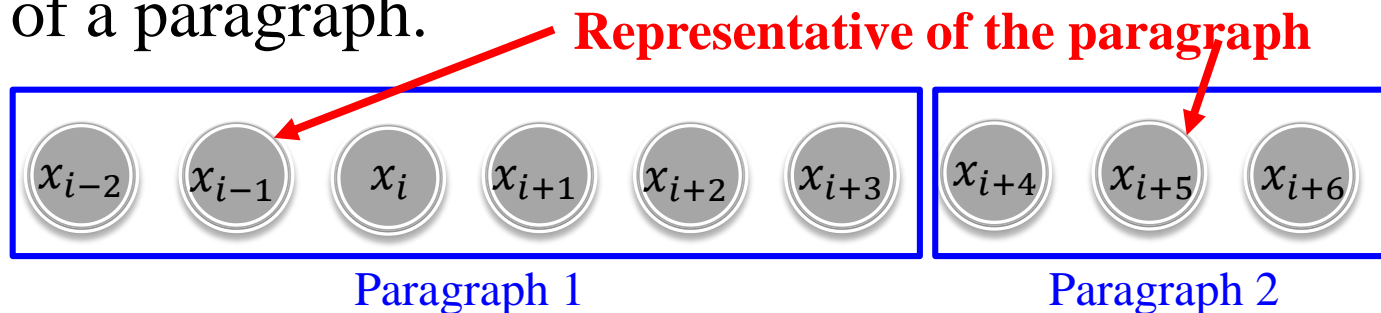
- Structure in Spoken Content

- Temporal structure helps summarization
 - ▣ Long summary: consecutive utterances in a paragraph are more likely to be

Important paragraph



- ▣ Short summary: one utterance is selected on behalf of a paragraph.



Summarization

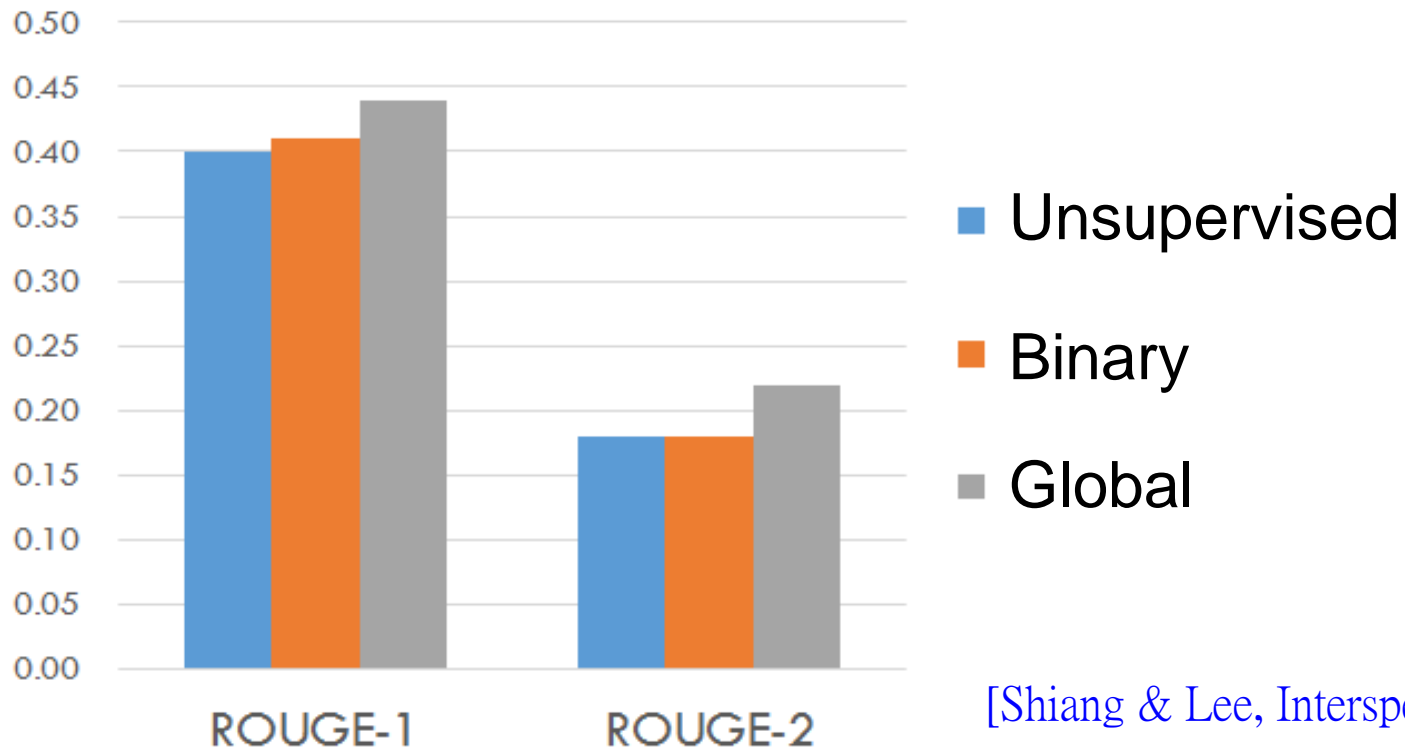
- Structure in Spoken Content

- Structure in text are clear
 - ▣ Paragraph boundaries are directly known
- **For spoken content, there is no obvious structure**
 - ▣ The structure can be considered as “hidden variables”
 - ▣ Jointly learning structure of spoken document and summarization [Shiang & Lee, Interspeech 13]

Summarization

- Experiments

- Evaluation Measure: ROUGE-1 and ROUGE-2
 - ▣ Larger scores means the machine-generated summaries is more similar to human-generated summaries.



[Shiang & Lee, Interspeech 13]

Key Term Extraction

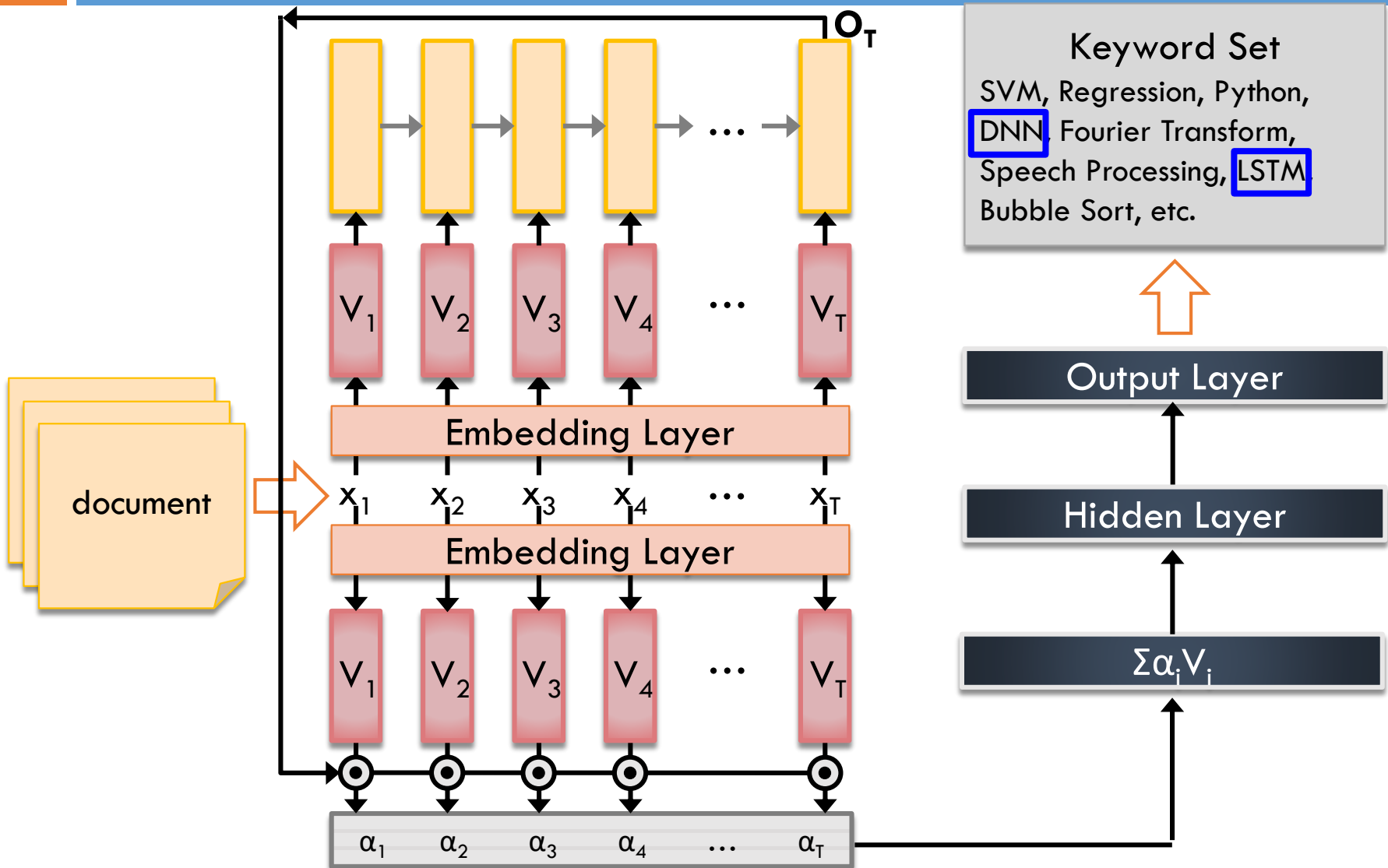
- TF-IDF is a good measure for identifying key terms
[E. D'Avanzo, DUC 04][Jiang, SIGIR 09]
- Feature parameters from latent topic models [Hazen, Interspeech 11] [Chen & Lee, SLT 10]
 - ▣ Key terms are usually focused on small number of topics
- Prosodic Features [Chen & Lee, ICASSP 12]
 - ▣ slightly lower speed, higher energy, wider pitch range
- Machine Learning methods
 - ▣ Input: a term, output: key term or not [Liu, SLT 08][Chen & Lee, SLT 10]
 - ▣ Input: a document, output: key terms in the document
[Kamal Sarkar, arXiv, 2010]

Key Term Extraction

– Deep Learning

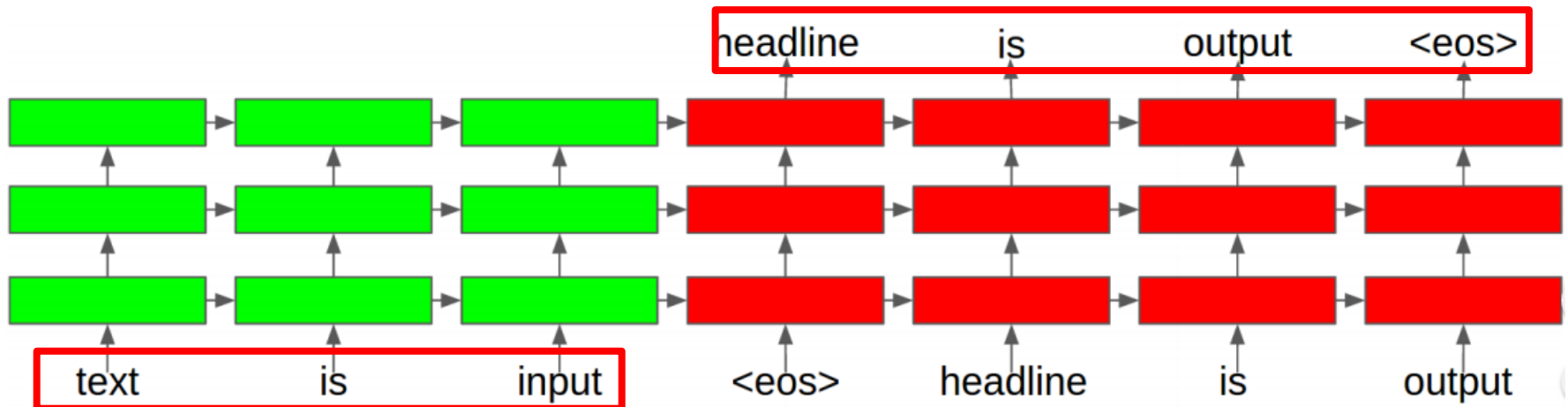
[Shen & Lee, Interspeech 16]

Poster, Sunday, 1:30 p.m., Dialogue Systems and Analysis of Dialogue



Title Generation

- Deep Learning based Approach [Alexander M Rush, EMNLP 15][Chopra, NAACL 16][Lopyrev, arXiv 2015][Shen, arXiv 2016]
 - ▣ Based on Sequence-to-sequence learning
 - ▣ Input: a document (word sequence), output: its title (shorter word sequence)



New Direction 5-3:

Speech Content is Difficult to Browse!

Organizing Retrieved Results

Introduction

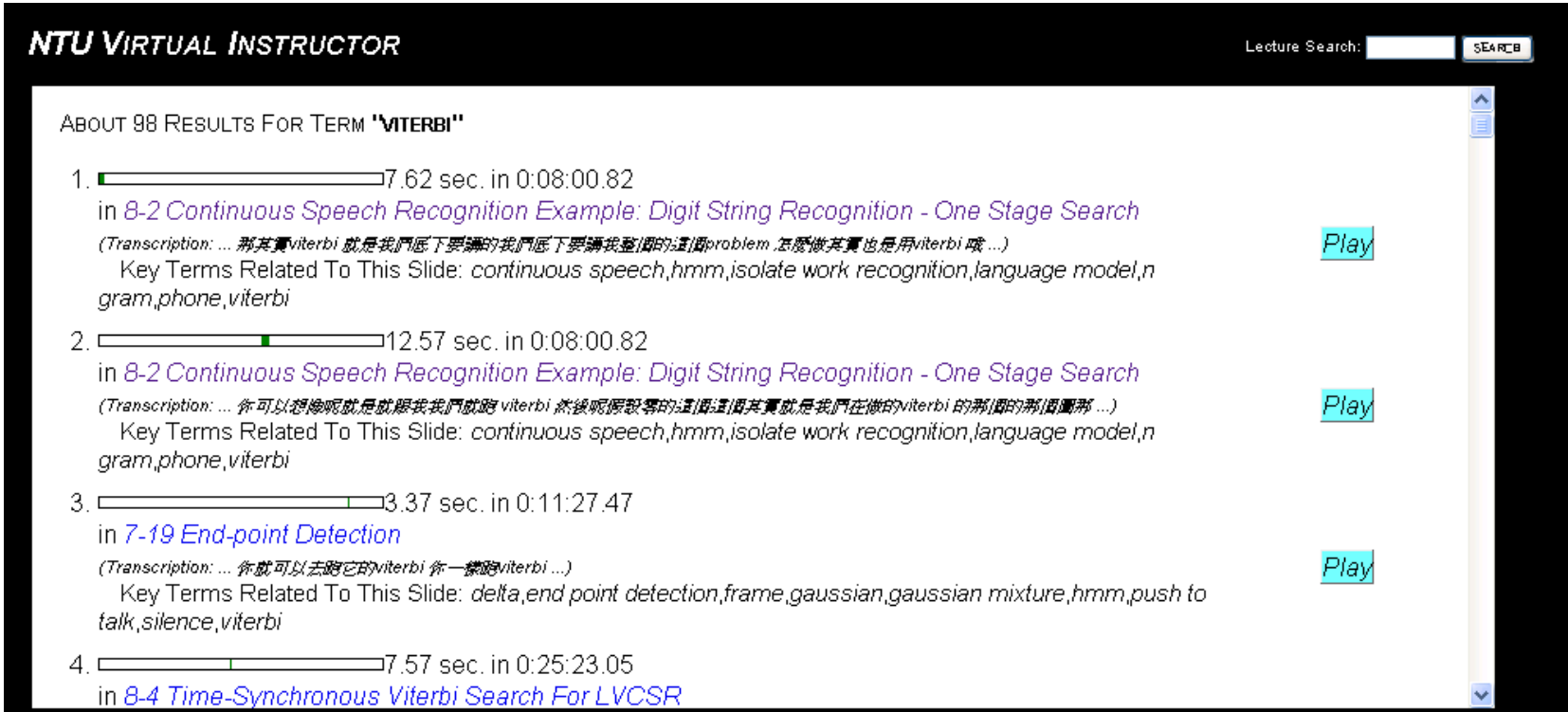
- Organizing the retrieval results to help users know what is retrieved
- Taking retrieving on-line lectures as example
 - ▣ Searching spoken lectures is a very good application for spoken content retrieval
 - ▣ The speech of the instructors conveys most knowledge in the lectures

Retrieving One Course

□ NTU Virtual Instructor

[Kong & Lee, ICASSP 09]

[Lee & Lee, IEEE/ACM T. ASL 14]

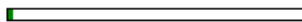





The screenshot shows the NTU Virtual Instructor interface. At the top left, it says "NTU VIRTUAL INSTRUCTOR". At the top right, there is a search bar labeled "Lecture Search:" with a "SEARCH" button. Below the search bar, it displays "ABOUT 98 RESULTS FOR TERM 'VITERBI'". There are four search results listed, each with a progress bar, a duration, a title, a transcription snippet, and key terms. To the right of each result is a "Play" button.

NTU VIRTUAL INSTRUCTOR

Lecture Search: SEARCH

ABOUT 98 RESULTS FOR TERM "VITERBI"

1.  7.62 sec. in 0:08:00.82
in [8-2 Continuous Speech Recognition Example: Digit String Recognition - One Stage Search](#)
(Transcription: ... 那其實viterbi 就是我們底下要講的我們底下要講我整個的這個problem 怎麼做其實也是用viterbi 嘍...)
Key Terms Related To This Slide: *continuous speech,hmm,isolate work recognition,language model,n gram,phone,viterbi* [Play](#)
2.  12.57 sec. in 0:08:00.82
in [8-2 Continuous Speech Recognition Example: Digit String Recognition - One Stage Search](#)
(Transcription: ... 你可以想像呢就是就跟我我們就聽 viterbi 然後呢像說零的這個這個其實就是我們在做的viterbi 的那個的那個...)
Key Terms Related To This Slide: *continuous speech,hmm,isolate work recognition,language model,n gram,phone,viterbi* [Play](#)
3.  3.37 sec. in 0:11:27.47
in [7-19 End-point Detection](#)
(Transcription: ... 你就可以去聽它的viterbi 你一樣聽viterbi...)
Key Terms Related To This Slide: *delta,end point detection,frame,gaussian,gaussian mixture,hmm,push to talk,silence,viterbi* [Play](#)
4.  7.57 sec. in 0:25:23.05
in [8-4 Time-Synchronous Viterbi Search For LVCSR](#)

Searching the course Digital Speech Processing of NTU

Massive Open On-line Courses (MOOCs)

- Enormous on-line courses



Today's Retrieval Techniques

coursera

Courses Specializations Institutions About | Hung-yi Lee ▾

language model

Global Partners (3) · US State Institutions (0)

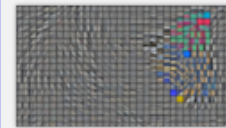
Sort by Starting soon ▾

- Always Open
- Starting Soon

Eligible For

- Specialization Certificates
- Verified Certificates
- All Partners
- Columbia University
- Stanford University
- University of Toronto
- All Languages
- English
- Arabic

Courses



University of Toronto
Neural Networks for Machine Learning
with Geoffrey Hinton

Oct 1st 2012
8 weeks long



Columbia University
Natural Language Processing
with Michael Collins

Feb 24th 2013
10 weeks long



Stanford University
Natural Language Processing
with Dan Jurafsky & Christopher Manning

There are
no open sessions.

A list of related courses

Today's Retrieval Techniques

coursera

Courses Specializations Institutions About | Hung-yi Lee ▾

Second-Order Markov Process

$$P(x_2, x_3, \dots, x_n | x_1) = \prod_{i=2}^n P(x_i | x_{i-1}, x_{i-2})$$

Typical Values of Perplexity

- Baseline baseline (1.0) for unigram language model: $P(x_i | x_{i-1}) = 1/|V|$
- A unigram model: $P(x_i | x_{i-1}) = 1/|V|$
- A bigram model: $P(x_i | x_{i-1}, x_{i-2}) = 1/|V|^2$
- A trigram model: $P(x_i | x_{i-1}, x_{i-2}, x_{i-3}) = 1/|V|^3$

Basic NLP Problem: Tagging

Courses

- 0 University of Toronto
Neural Networks for Machine Learning with Geoffrey Hinton
- 0 Columbia University
Natural Language Processing with Michael Collins
- 3 Stanford University
Natural Language Processing with Dan Jurafsky & Christopher Manning

Neural Networks for Machine Learning

Lecture 1a
Why do we need machine learning?

Geoffrey Hinton with Nick Srivastava, Kevin Swersky

Ways to learn

- A large number of different methods
- Weight sharing
- Early stopping
- Model averaging
- Bayesian fitting of parameters
- Dropout
- Many of these methods

Introduction to NLP

What is Natural Language Processing?

More sophisticated decision tree features

- Case of word with "": Upper, Lower, Number
- Case of word after "": Upper, Lower, Cap, Number
- Numeric features
- Length of word with "":
- Probability of word with "": occurs at end of s
- Probability of word after "": occurs at beginning of s

Applying Multinomial Naive Bayes Classifiers to Text Classification

positions ← all word positions in test document

$$c_{\text{top}} = \underset{c \in C}{\text{argmax}} P(c) \prod_{i=1}^n P(x_i | c)$$

$$c_i = \frac{P(c_i)}{P(c_i) + \prod_{j \neq i} P(c_j)}$$

More is less

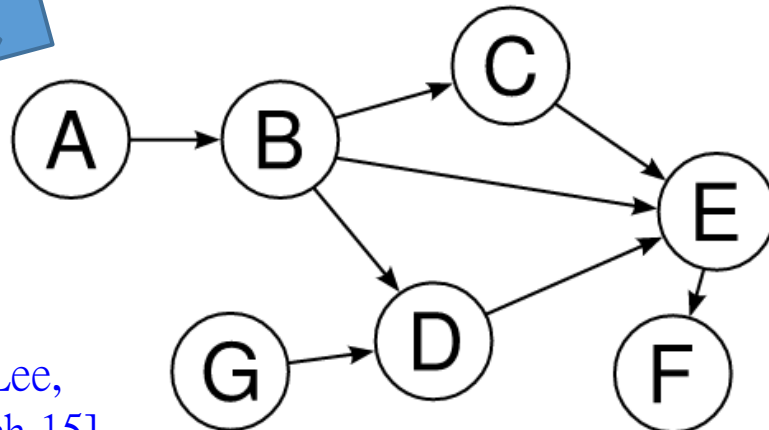
- Given all the related lectures from different courses



Which lecture should I go first?



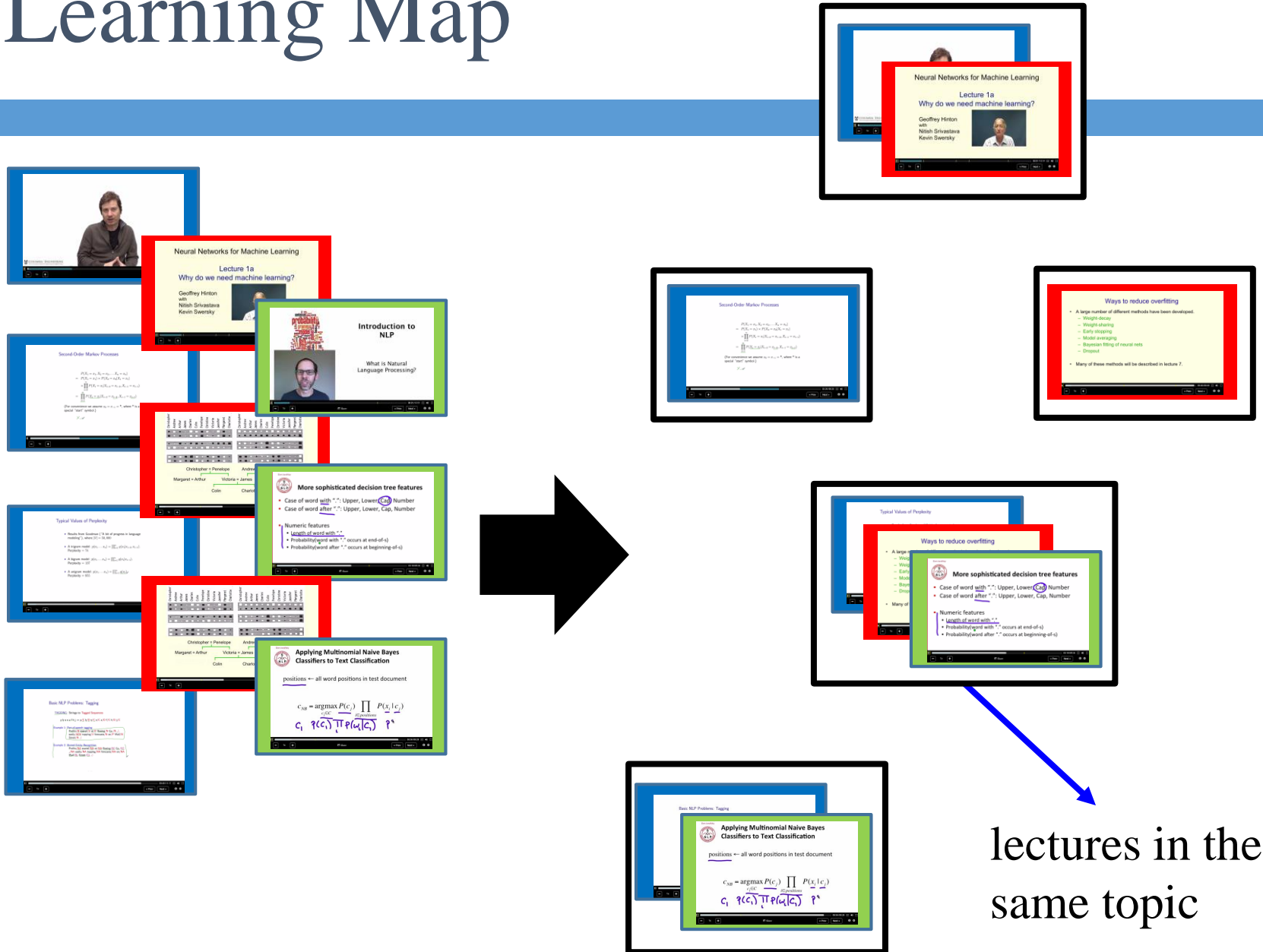
learner



Learning Map

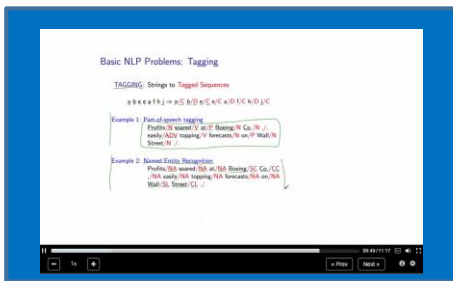
- Nodes: lectures in the same topics
- Edges: suggested learning order

Learning Map

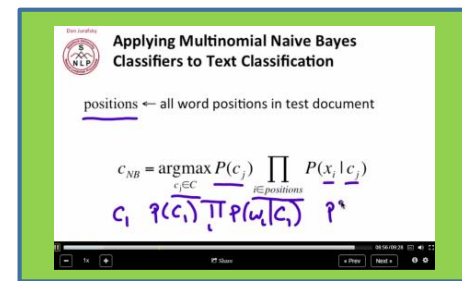


Lectures in the same topic

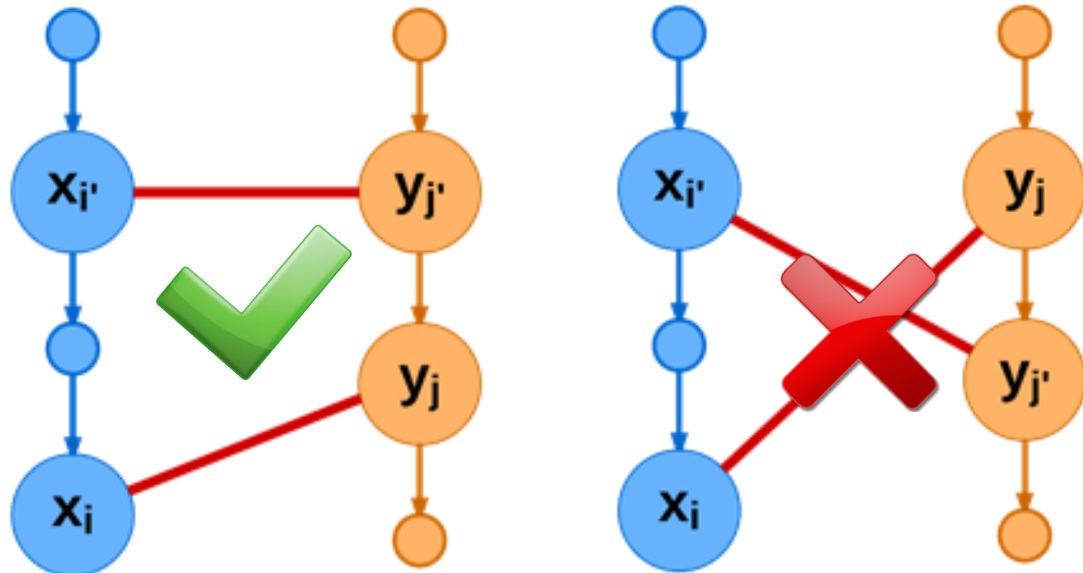
“Local” Information:



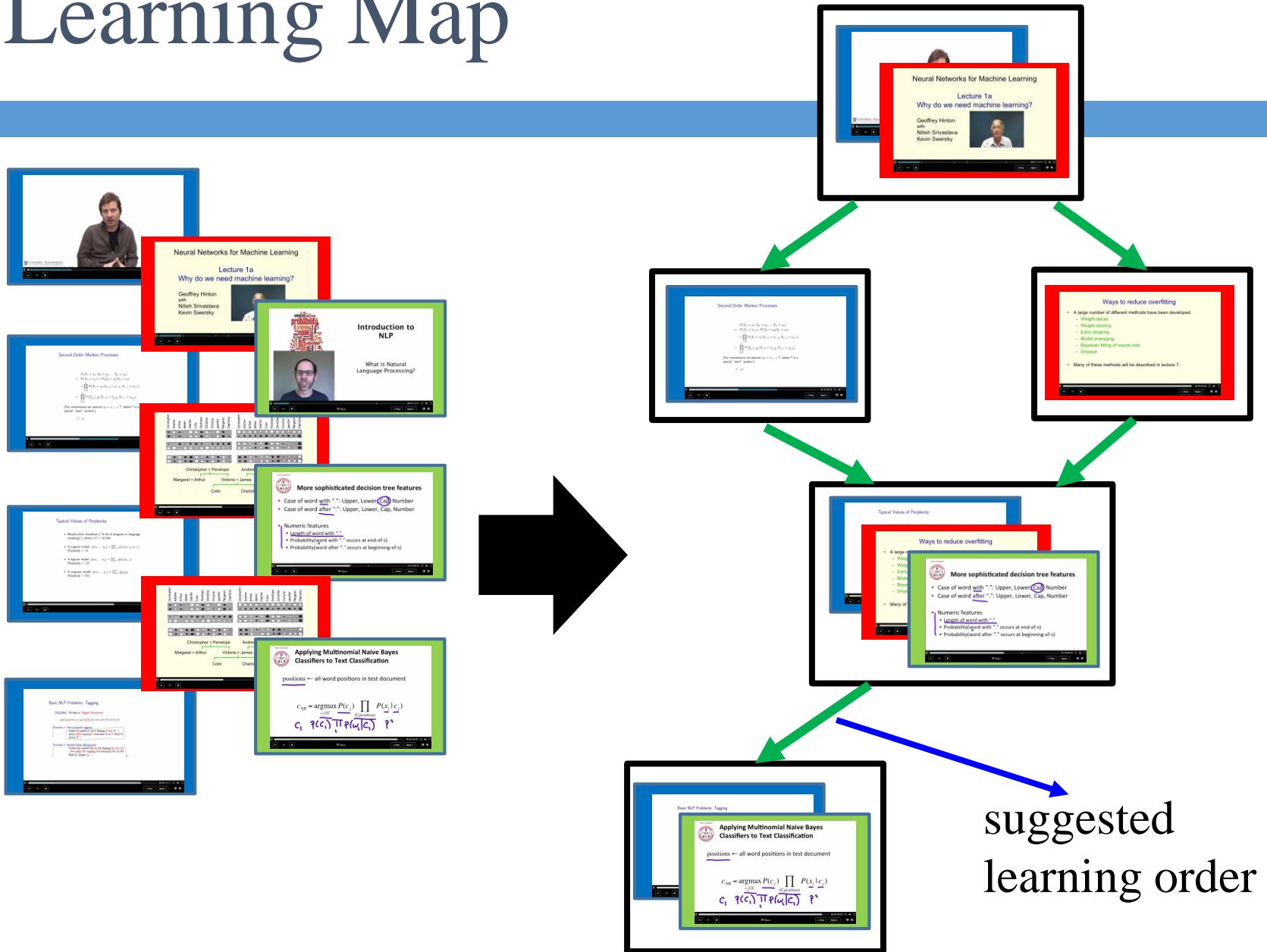
Similarity?



“Global” Information:



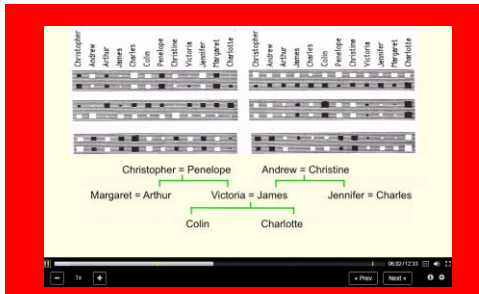
Learning Map



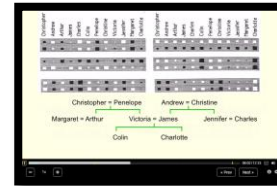
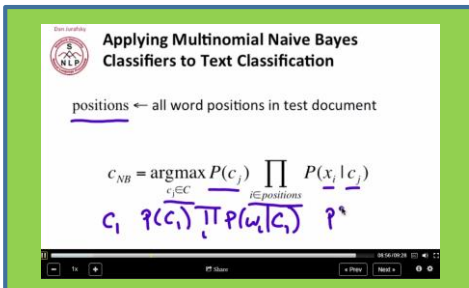
suggested learning order

Prerequisite

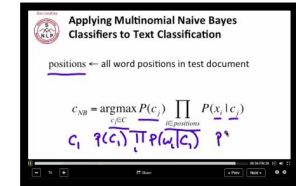
Lectures in different courses



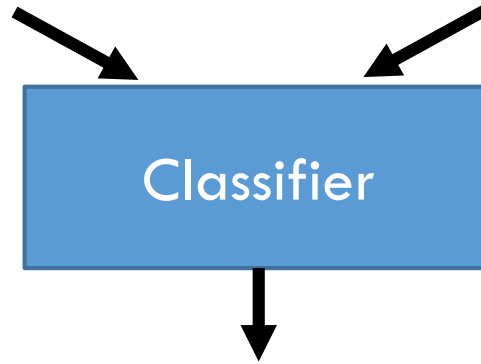
Prerequisite?



Content of Lecture A

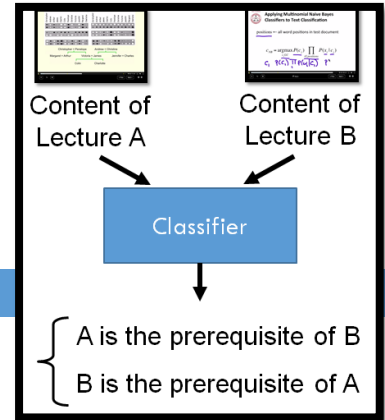


Content of Lecture B

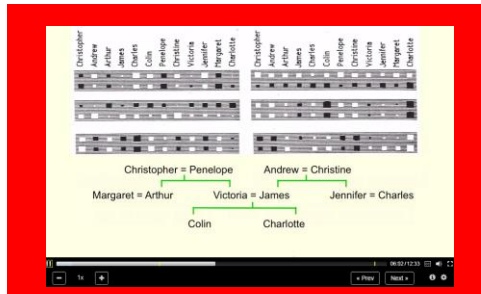


A is the prerequisite of B
B is the prerequisite of A

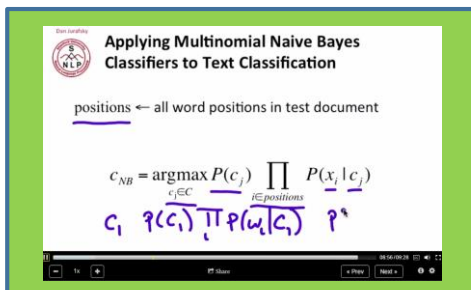
Prerequisite



Lectures in different courses

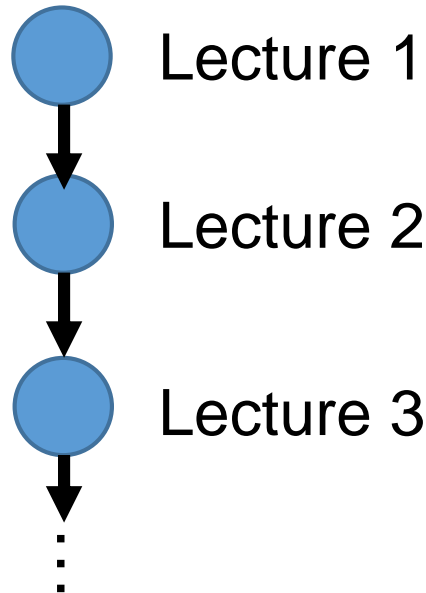


Prerequisite?



The existing courses on-line can be the training data

An existing course



Lecture 1 is a prerequisite of lecture 2
 Lecture 2 is a prerequisite of lecture 3

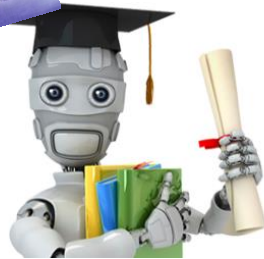
Training examples

Demo



Vision: Personalized Courses

on-line learning
material



- I want to learn “deep learning”.
- I am a graduate student of computer science.
- I can spend 6 hours.



Learner

I open a course for you.

- With MOOCs and Spoken Language Processing techniques
 - ▣ It is possible to have a personalized course for each learning need.

New Direction 5-4:

Speech Content is Difficult to Browse!

Spoken Question Answering

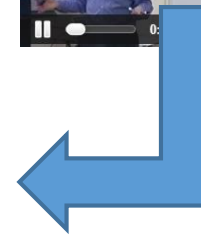
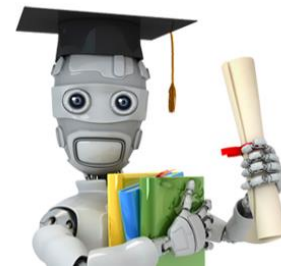
Spoken Question Answering

What is a possible origin of Venus' clouds?

Venus' clouds

Taking some time to find the answer

Lectures
about Venus



Spoken Content Retrieval

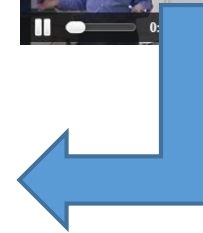
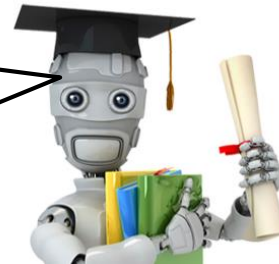
Spoken Question Answering



What is a possible origin of Venus' clouds?



Gases released as a result of volcanic activity



Spoken Question Answering: Machine answers questions based on the information in spoken content

Spoken Question Answering

- Question Answering in Speech Transcripts (QAST) has been a well-known evaluation program of spoken question answering.
 - ▣ 2007, 2008, 2009

Reference: **6 Spoken Question Answering** (*Sophie Rosset, Olivier Galibert and Lori Lamel*). G. Tur and R. DeMori, Spoken Language Understanding: Systems for Extracting Semantic Information from Speech.

- Focused on factoid questions in the previous study
 - ▣ E.g. “What is name of the highest mountain in Taiwan?”.
- To answer more difficult questions, machine has to understand questions and spoken documents.
 - ▣ How good can it achieve?

New task for Machine Comprehension of Spoken Content

□ TOEFL Listening Comprehension Test by Machine

[Tseng & Lee, Interspeech 16]

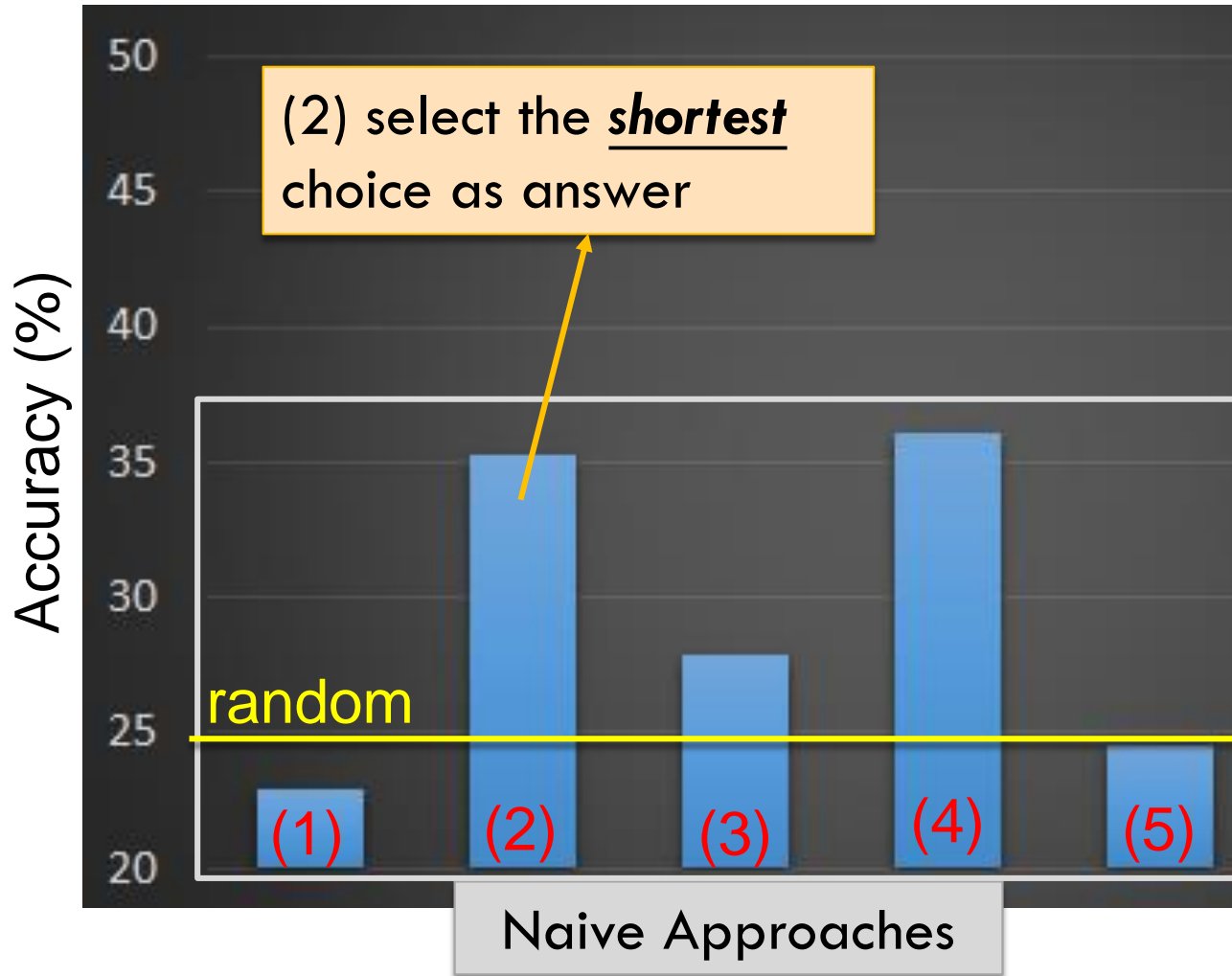
Audio Story:  (The original story is 5 min long.)

Question: “ What is a possible origin of Venus’ clouds? ”

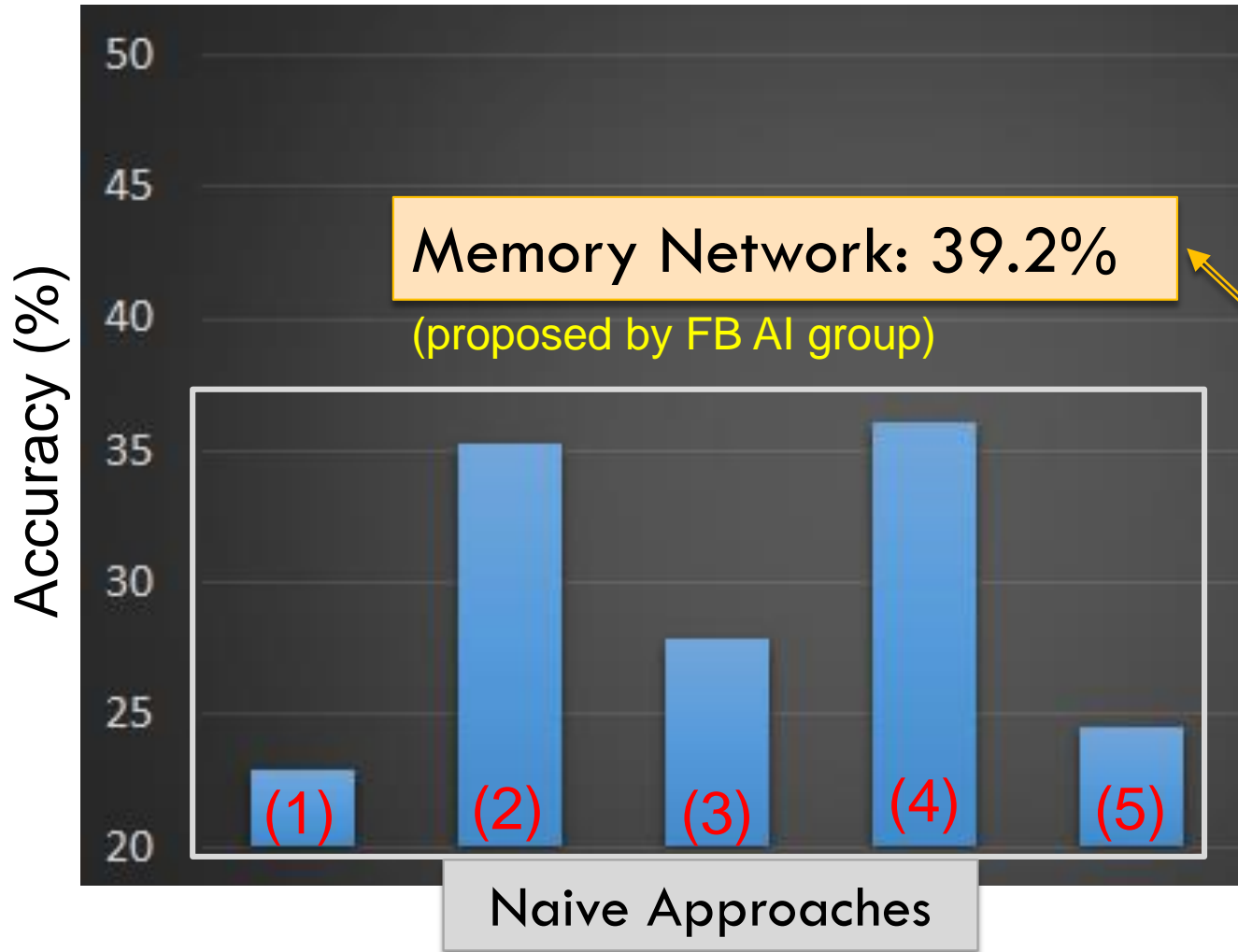
Choices:

- (A) gases released as a result of volcanic activity
- (B) chemical reactions caused by high surface temperatures
- (C) bursts of radio energy from the plane's surface
- (D) strong winds that blow dust into the atmosphere

Simple Baselines

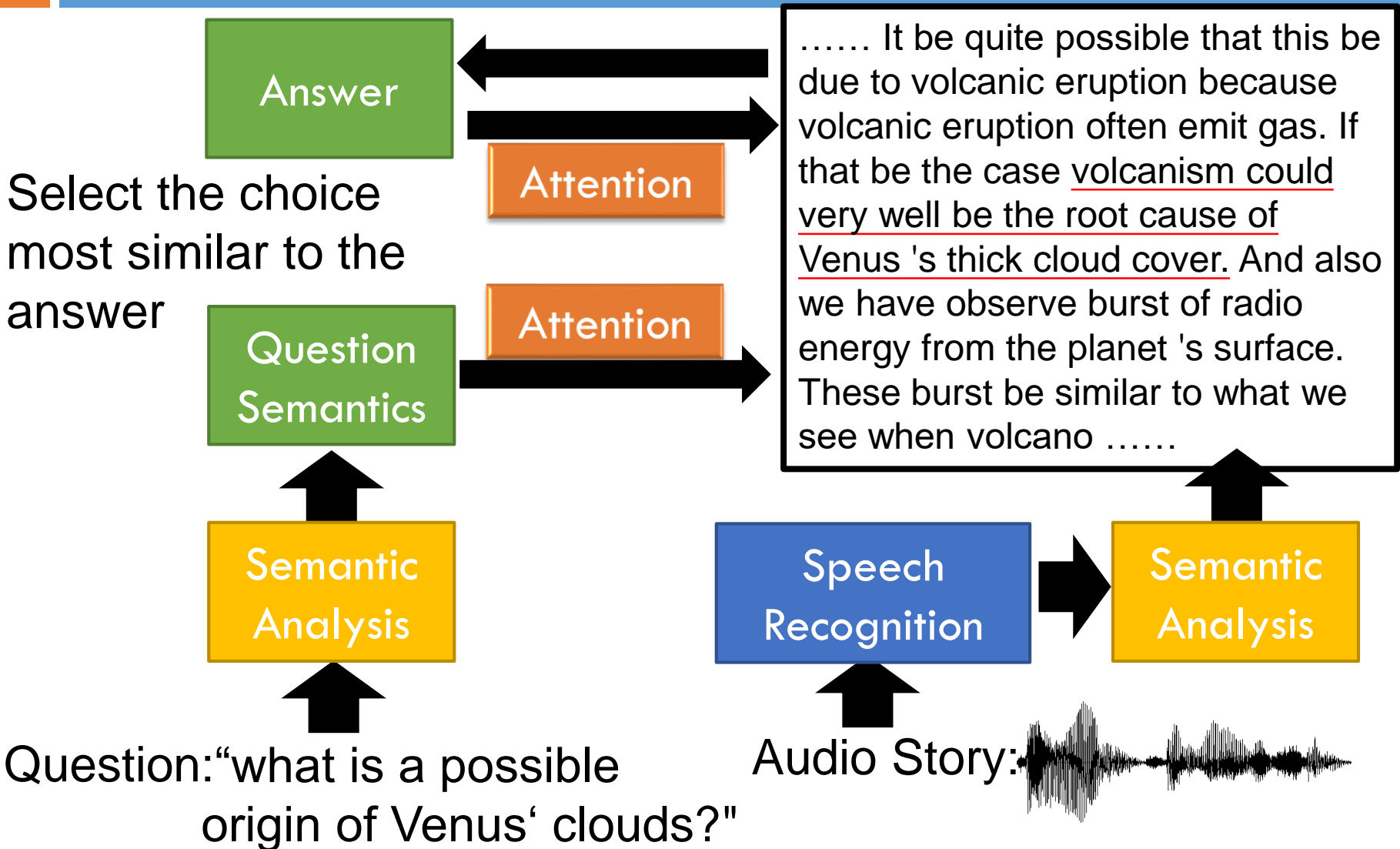


Results

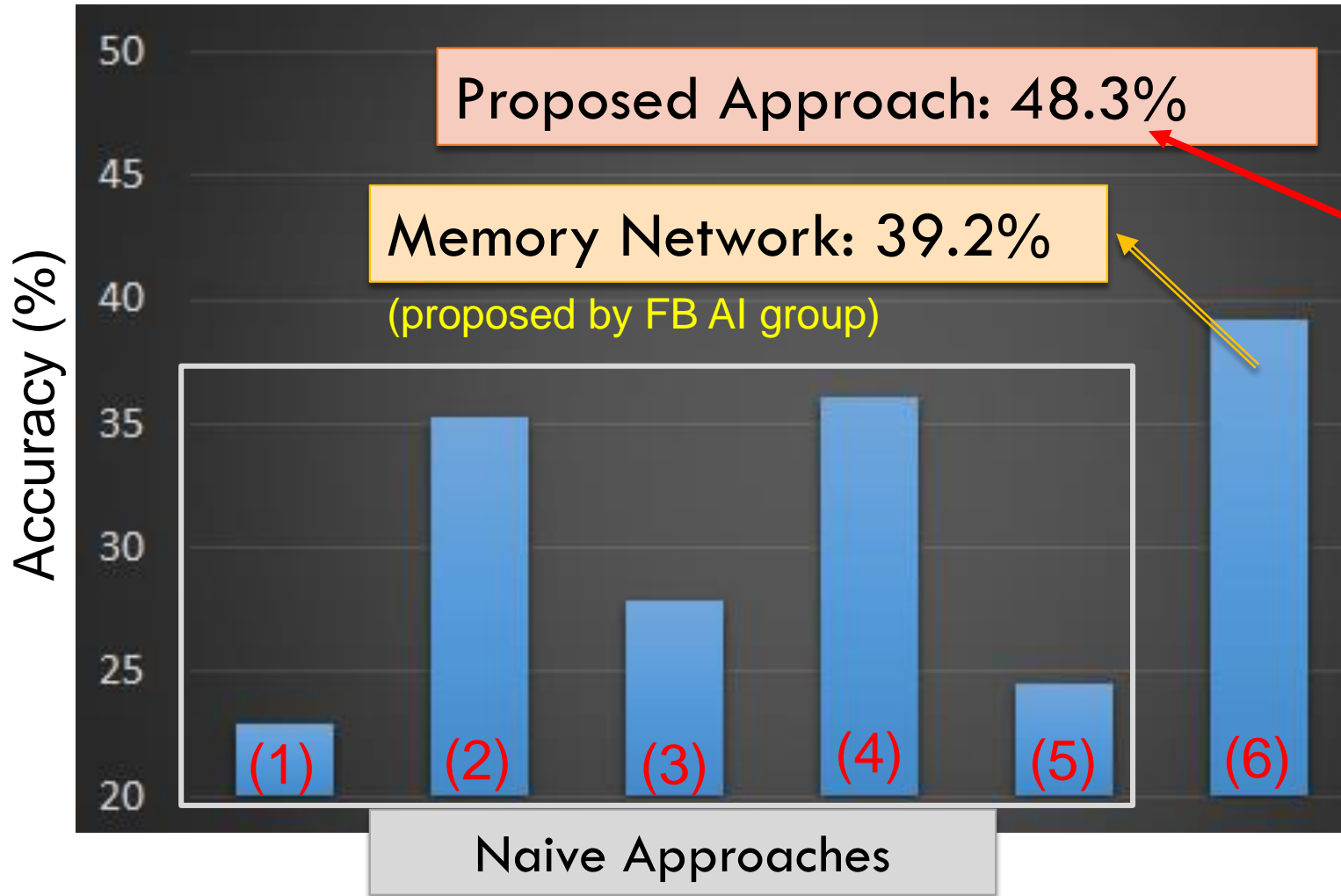


Model Architecture

The model is learned end-to-end.



Results



Concluding Remarks

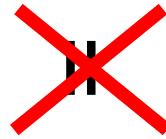


Conclusion Remarks

- New research directions for spoken content retrieval
 - ▣ Modified ASR for Retrieval Purposes
 - ▣ Incorporating Those Information Lost in ASR
 - ▣ No Speech Recognition!
 - ▣ Special Semantic Retrieval Techniques for Spoken Content
 - ▣ Spoken Content is Difficult to Browse!

Take-Home Message

Spoken Content Retrieval



Speech Recognition

+

Text Retrieval

Spoken Content Retrieval



300 hrs multimedia is
uploaded per minute.
(2015.01)




1874 courses on coursera
(2016.04)

- Nobody is able to go through the data.
- In these multimedia, the spoken part carries very important information about the content
- Spoken content retrieval: Machine listens to the data, and extract the desired information for each individual user.
 - Just as Google does on text data

Overview Paper

- Lin-shan Lee, James Glass, Hung-yi Lee, Chun-an Chan, "Spoken Content Retrieval —Beyond Cascading Speech Recognition with Text Retrieval," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol.23, no.9, pp.1389-1420, Sept. 2015
- <http://speech.ee.ntu.edu.tw/~tlkagk/paper/Overview.pdf>
- This tutorial includes updated information after this paper is published.



Thank You for Your Attention