

Posterior-based Approach

Dr. Mathew Magimai Doss

Idiap Research Institute, Martigny, Switzerland

Posterior-based ASR Theory (1)

$$W^* = \arg \max_{W_k} P(W_k | X)$$

$$\begin{aligned} P(W_k | X) &= \sum_Q P(W_K, Q | X) & X &= \{\mathbf{x}_1, \dots, \mathbf{x}_m, \dots, \mathbf{x}_M\} \\ &= \sum_Q P(W_k | X, Q) \cdot P(Q | X) & Q &= \{q_1, \dots, q_{m-1}, \dots, q_M\} \\ &\simeq \sum_Q P(Q | X) \cdot P(W_K | Q) \end{aligned}$$

Posterior-based ASR Theory (2)

$$P(Q|X) = \prod_{m=1}^M P(q_m|X, Q_1^{m-1}) \quad \text{applying Gibbs rule} \quad Q_1^{m-1} = \{q_1, \dots, q_{m-1}\}$$

$$\begin{aligned} P(W_k|Q) &= \frac{P(Q|W_k) \cdot P(W_k)}{P(Q)} \quad \text{applying Bayes' rule} \\ &= \frac{(\prod_{m=1}^M P(q_m|Q_1^{m-1}, W_k) \cdot P(W_k))}{\prod_{m=1}^M P(q_m|Q_1^{m-1})} \quad \text{applying Gibbs rule} \end{aligned}$$

$$P(W_k|X) \simeq \sum_Q \left(\prod_{m=1}^M P(q_m|X, Q_1^{m-1}) \right) \cdot \left(\frac{(\prod_{m=1}^M P(q_m|Q_1^{m-1}, W_k) \cdot P(W_k))}{\prod_{m=1}^M P(q_m|Q_1^{m-1})} \right)$$

$$P(W_k|X) \approx \max_Q \left(\prod_{m=1}^M P(q_m|X, Q_1^{m-1}) \right) \cdot \left(\frac{(\prod_{m=1}^M P(q_m|Q_1^{m-1}, W_k) \cdot P(W_k))}{\prod_{m=1}^M P(q_m|Q_1^{m-1})} \right) \quad \text{Viterbi approx.}$$

Posterior-based ASR Theory (3)

$$P(W_k|X) \simeq \sum_Q \left(\prod_{m=1}^M P(q_m|X, q_{m-1}) \cdot \left(\frac{\left(\prod_{m=1}^M P(q_m|q_{m-1}, W_k) \cdot P(W_k) \right)}{\prod_{m=1}^M P(q_m|q_{m-1})} \right) \right)$$

$$P(W_k|X) \approx \max_Q \left(\prod_{m=1}^M P(q_m|X, q_{m-1}) \cdot \left(\frac{\left(\prod_{m=1}^M P(q_m|q_{m-1}, W_k) \cdot P(W_k) \right)}{\prod_{m=1}^M P(q_m|q_{m-1})} \right) \right)$$

$P(q_m|X, q_{m-1})$ = Conditional prior probabilities

$P(q_m|q_{m-1}, W_k)$ = similar to HMM transition prob.

$P(q_m|q_{m-1})$ = Training data priors (time-invariant HMM)

Acoustic modeling, language modeling and decoder are not independent any more

Hybrid HMM/ANN ASR

$$\begin{aligned} P(W_k|X) &\approx \max_Q \left(\prod_{m=1}^M P(q_m|X, q_{m-1}) \cdot \left(\frac{\left(\prod_{m=1}^M P(q_m|q_{m-1}, W_k) \cdot P(W_k) \right)}{\prod_{m=1}^M P(q_m|q_{m-1})} \right) \right. \\ &\approx \max_Q \left(\prod_{m=1}^M P(q_m|X_{m-c}^{m+c}) \right) \cdot \left(\frac{\left(\prod_{m=1}^M P(q_m|q_{m-1}, W_k) \cdot P(W_k) \right)}{\prod_{m=1}^M P(q_m)} \right) \\ &\approx \max_Q \left(\prod_{m=1}^M \frac{P(q_m|X_{m-c}^{m+c})}{P(q_m)} \cdot P(q_m|q_{m-1}, W_k) \right) \cdot P(W_k) \end{aligned}$$

$$X_{m-c}^{m+c} = \{\mathbf{x}_{m-c} \cdot \mathbf{x}_m, \dots, \mathbf{x}_{m+c}\} \quad \text{Limited context}$$

$$\frac{P(q_m|X_{m-c}^{m+c})}{P(q_m)} = \text{emission scaled-likelihood}$$

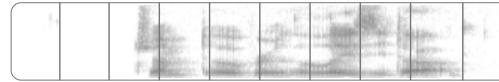
$$P(q_m|q_{m-1}, W_k) = \text{HMM transition prob.}$$

$$P(q_m | X)$$

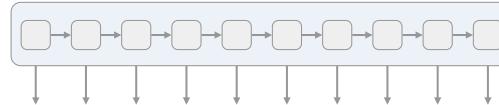
CTC

- Drop dependency on previous states
- Model full acoustic context using transformers
- Frame shift: 40 ms (25 frames per second)
 - With 10 ms frame shift the approach can lead to inferior performance
- Introduce empty symbols to emulate string matching at the output

Connectionist Temporal Classification



We start with an input sequence, like a spectrogram of audio.



The input is fed into an RNN, for example.

h	h	h	h	h	h	h	h	h	h
e	e	e	e	e	e	e	e	e	e
l	l	l	l	l	l	l	l	l	l
o	o	o	o	o	o	o	o	o	o
€	€	€	€	€	€	€	€	€	€

The network gives $p_t(a | X)$, a distribution over the outputs $\{h, e, l, o, \epsilon\}$ for each input step.

h	e	€	l	l	€	l	l	o	o
h	h	e	l	l	€	€	l	€	o
€	e	€	l	l	€	€	l	o	o

With the per time-step output distribution, we compute the probability of different sequences

h	e	l	l	o
e	l	l	o	
h	e	l	o	

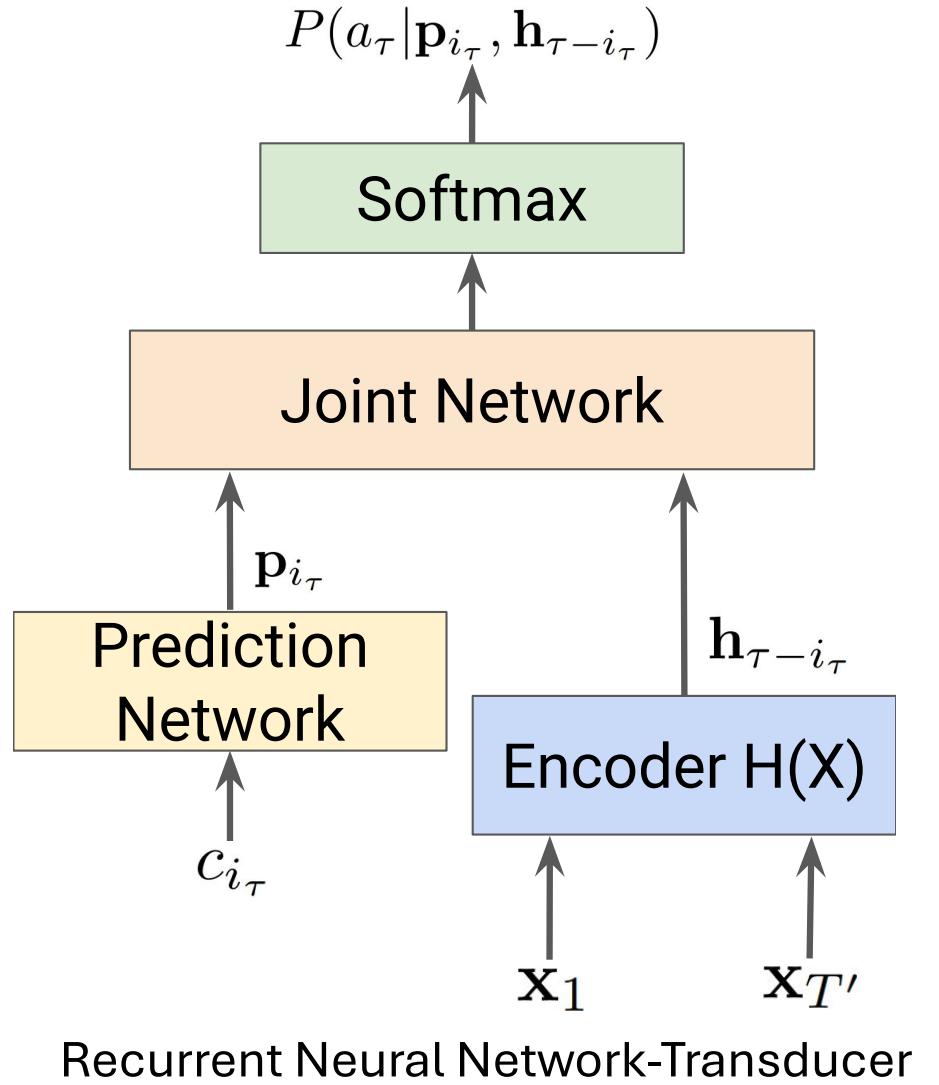
By marginalizing over alignments, we get a distribution over outputs.

Source: Hannun, "Sequence Modeling with CTC", Distill, 2017.

$$\left(\prod_{m=1}^M P(q_m | X, Q_1^{m-1}) \right)$$

RNN-T

- Model dependency on previous states (Prediction Network)
- Model full acoustic context using transformers (Encoder)
- Frame shift: 40 ms (25 frames per second)
- Introduce empty symbols to emulate string matching at the output



Sequence recognition

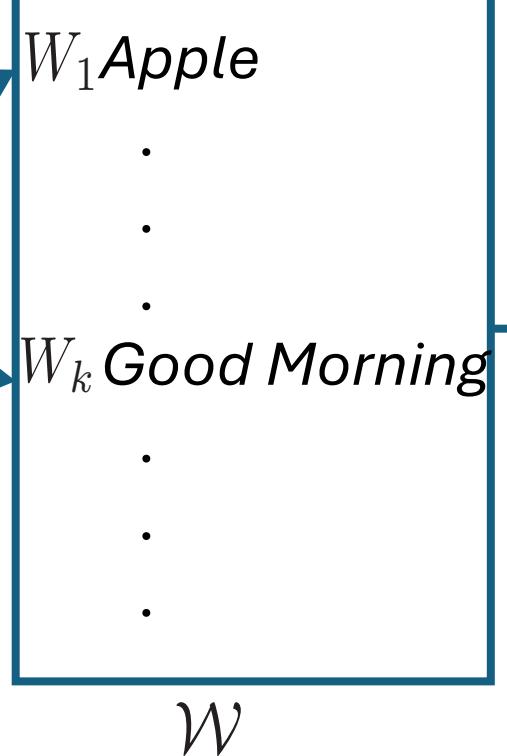
Automatic speech recognition (ASR)

Input



S

Set of hypotheses



Output

Best
Matching
Hypothesis
 \hat{W}

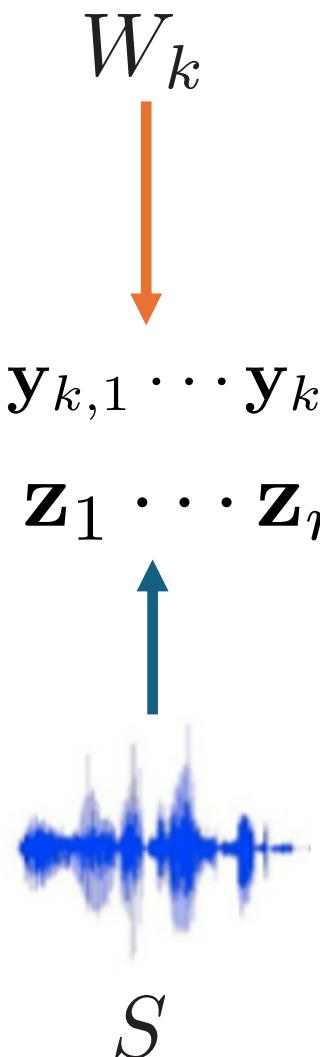
- **How to model the set of word hypothesis \mathcal{W} ?**
 - Sequence generation
Rule-based, discrete Markov model (DMM), neural LM, LLM
- **How to match W_k and S ?**
 - Sequence matching
- **How to search?**
 - Greedy search, beam search, stack decoding

$$\hat{W} = \arg \max_{W_k \in \mathcal{W}} \text{Match}(W_k, S)$$

Sequence matching: Four fundamental questions

Core idea

1. Map W_k and S to a shared latent symbol space.
 $Y_k = \mathbf{y}_{k,1} \cdots \mathbf{y}_{k,n} \cdots \mathbf{y}_{k,N}$
2. Match the resulting two latent symbol sequences Y_k and Z .
 $Z = \mathbf{z}_1 \cdots \mathbf{z}_m \cdots \mathbf{z}_M$



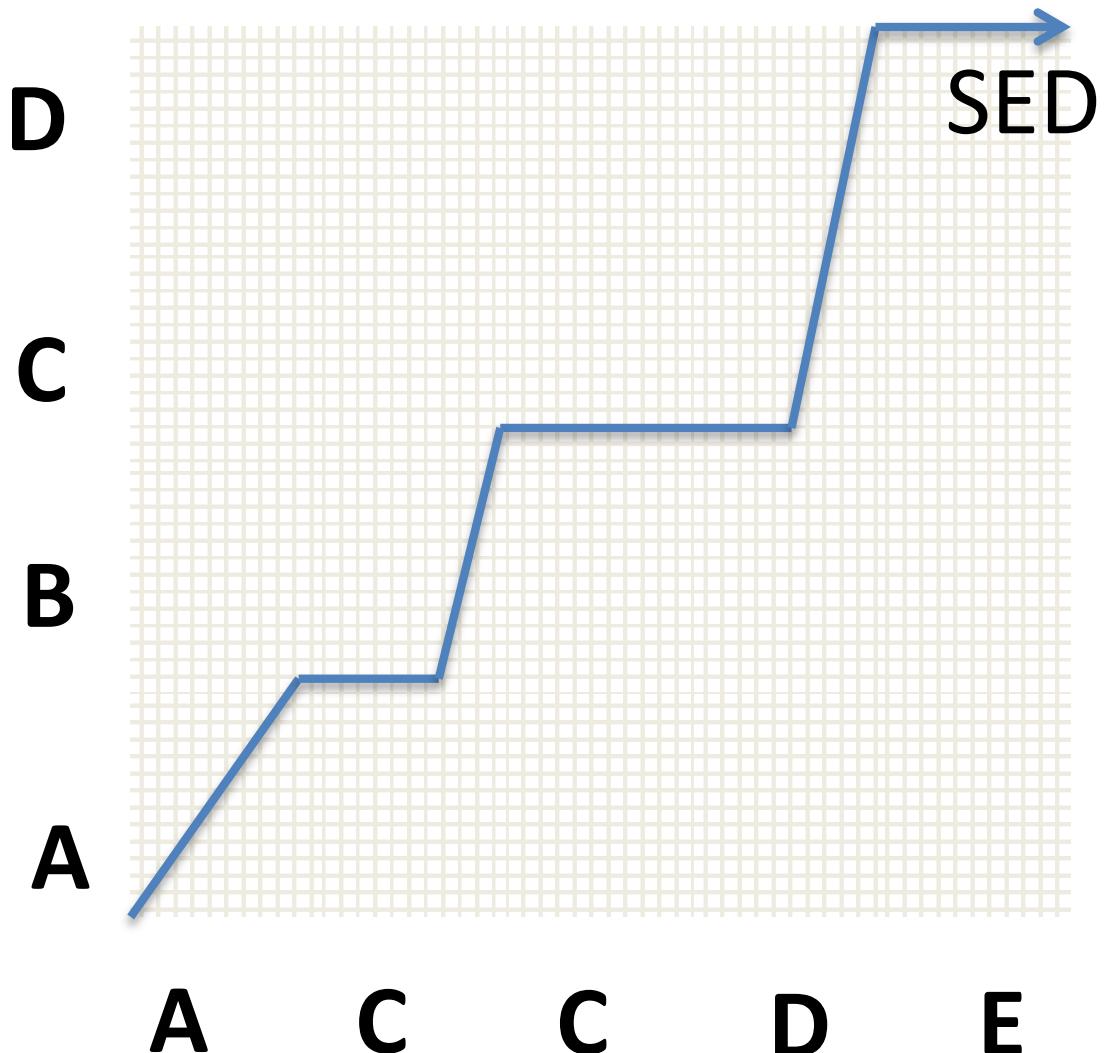
Q1: What is the shared latent symbol set $\{a^d\}_{d=1}^D$?

Q2: How to map S to latent symbol sequence Z ?

Q3: How to map W_k to latent symbol sequence Y_k ?

Q4: How to match Z and Y_k to estimate $\text{Match}(W_k, S)$?

String matching Q4



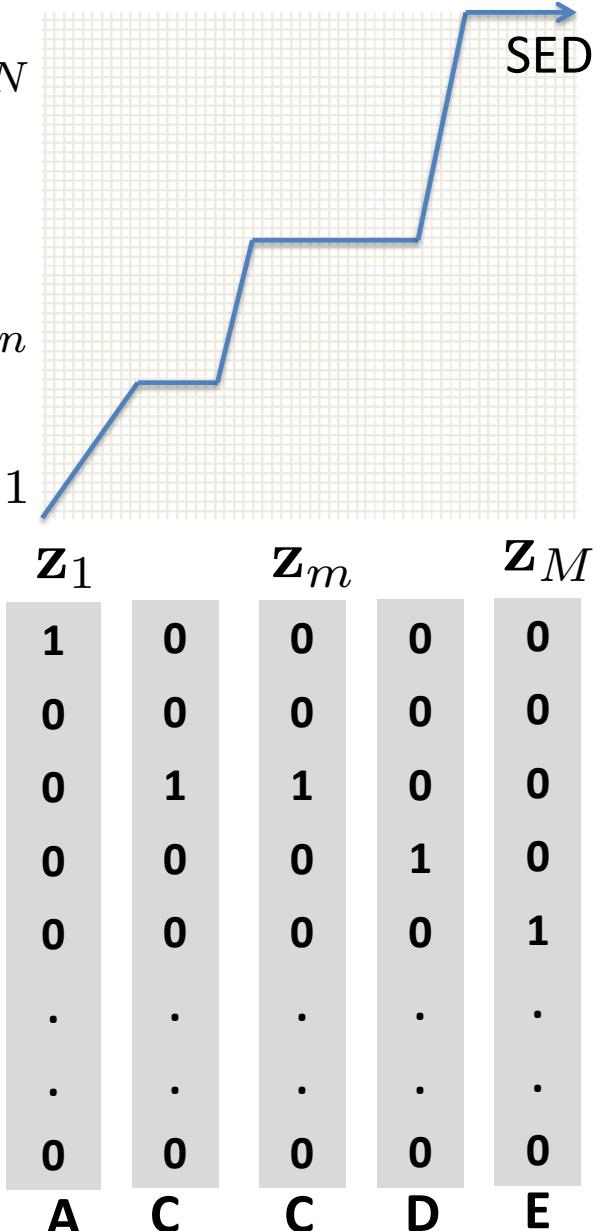
Initial conditions
for n from 1 to N
for m from 1 to M
if $\text{str}[m] = \text{str}[n]$
 $d[m,n] = 0$ // **Local score**
else
 $d[m,n] = 1$
done
 $D[m, n] = \min(\text{minimum}($
 $D[m-1, n] + 1, \text{// deletion}$
 $D[m, n-1] + 1, \text{// insertion}$
 $D[m-1, n-1] + d[m,n] \text{// substitution}$
 $) \text{// Global score}$

done
done

String edit distance SED = $D[M,N]$

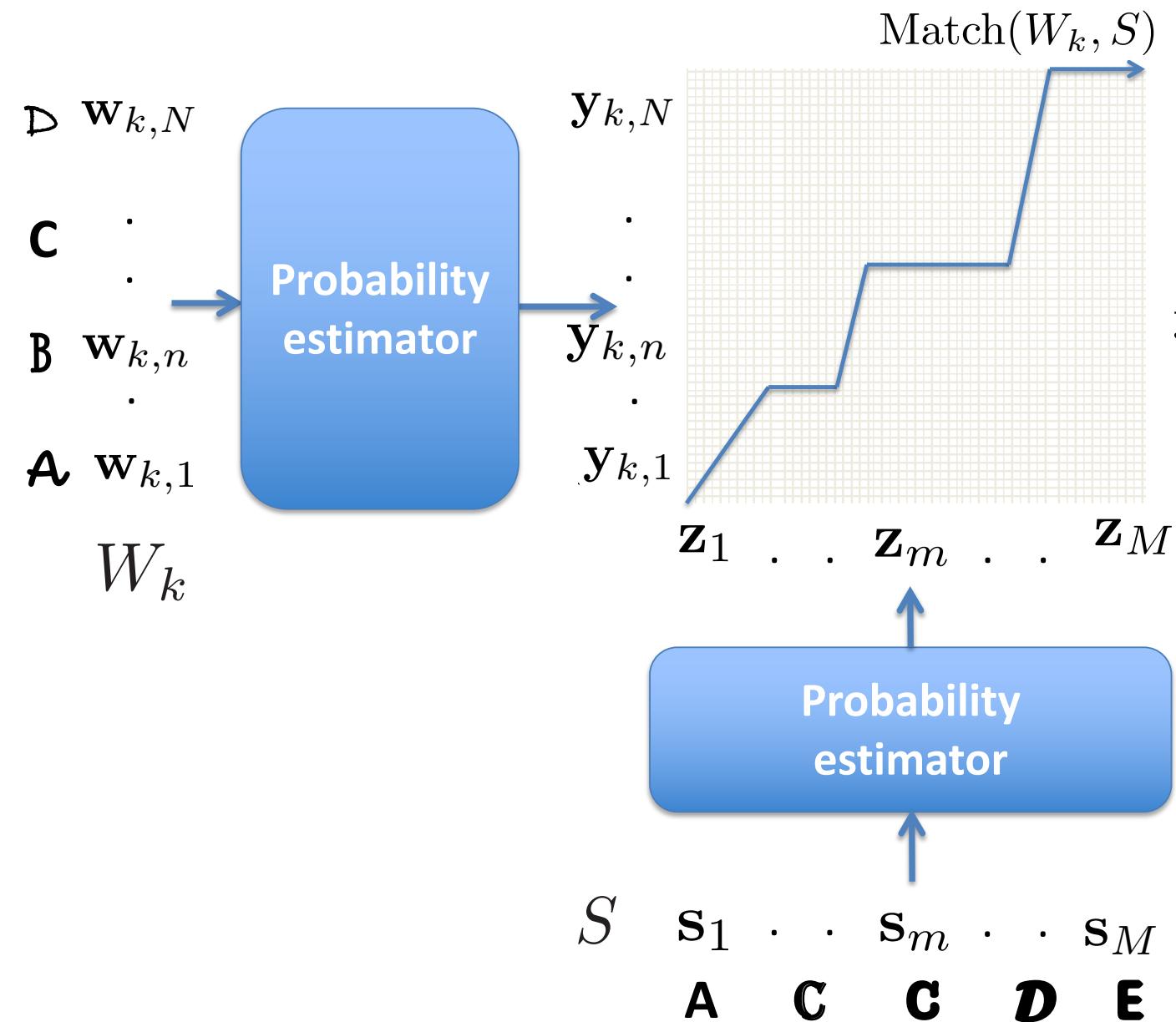
Deterministic symbol sequence matching

D 0 0 0 1 0 . . 0 $y_{k,N}$
C 0 0 1 0 0 . . 0
B 0 1 0 0 0 . . 0 $y_{k,n}$
A 1 0 0 0 0 . . 0 $y_{k,1}$



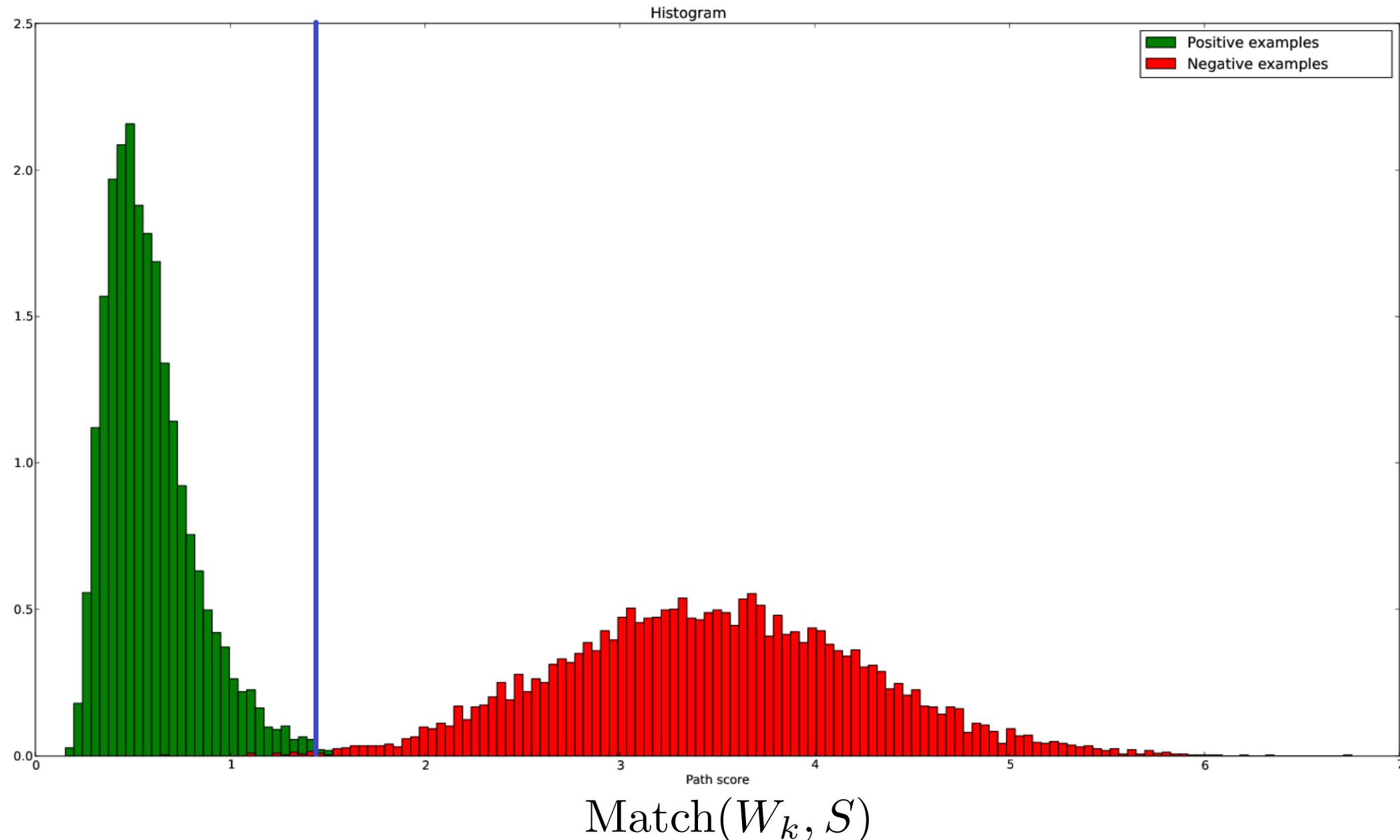
- Set of symbols
 $\{a^d\}_{d=1}^D = \{A, \dots, Z\}$
 - Deterministic symbols
 $y_{k,n} = \delta^d \quad z_m = \delta^{d'}$
 - Local score $d[m,n]$
if($KL(y_{k,n}, z_m) \leq \Delta$) $\Delta = 0$
 $d[m, n] = 0$
else
 $d[m, n] = 1$
- Kullback-Leibler (KL) divergence**
- $$KL(y_{k,n}, z_m) = \sum_{d=1}^D y_{k,n}^d \cdot \log\left(\frac{y_{k,n}^d}{z_m^d}\right)$$

Probabilistic symbol sequence matching

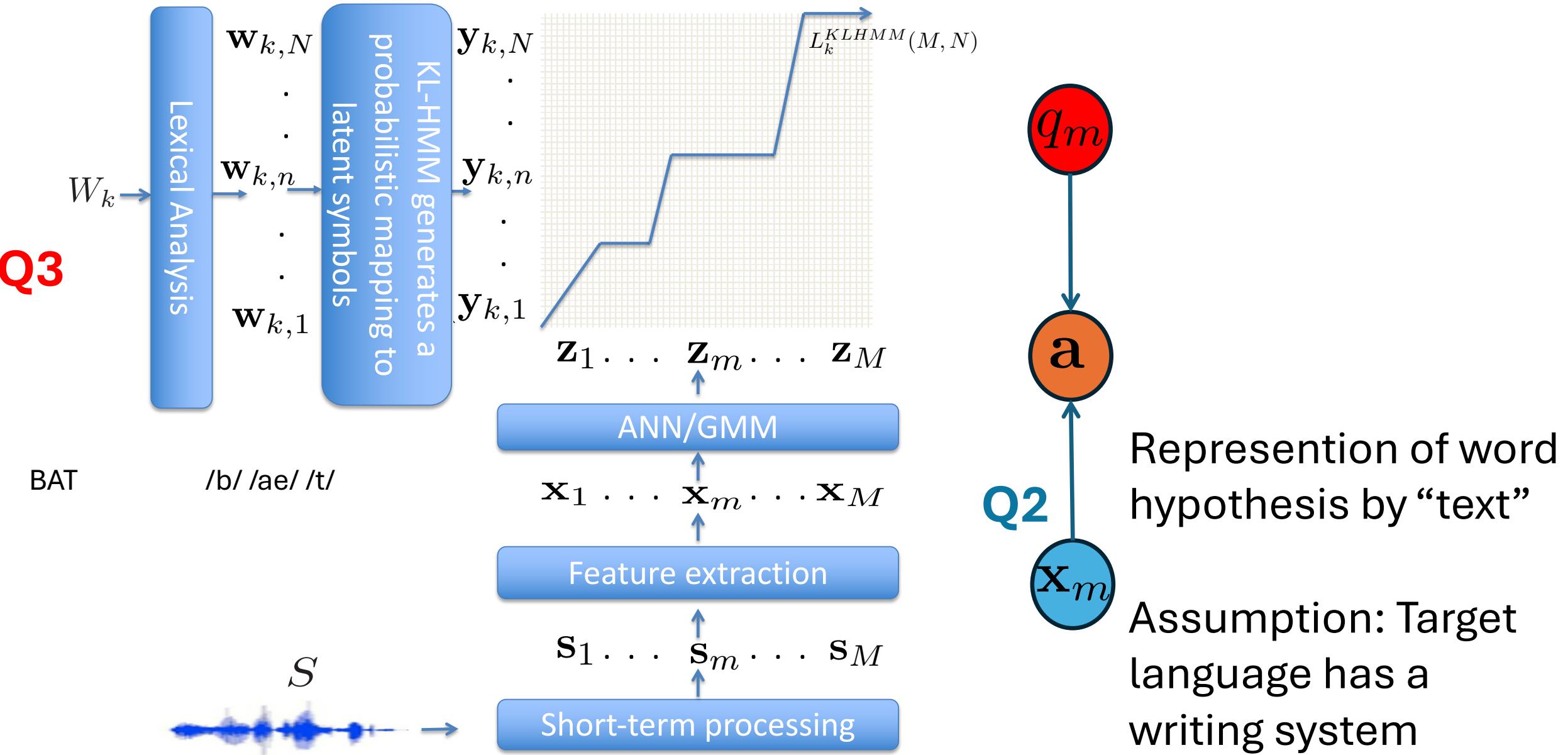


- Set of symbols
 $\{a^d\}_{d=1}^D = \{A, \dots, Z\}$
- **Posterior features** (Prob. symbols)
 $\mathbf{y}_{k,n} = [P(a^1|w_{k,1}) \dots P(a^d|w_{k,n}) \dots P(a^D|w_{k,N})]^T$
 $\mathbf{z}_m = [P(a^1|s_1) \dots P(a^d|s_m) \dots P(a^D|s_M)]^T$
- Local score $d[m,n]$
 $d[m, n] = \text{KL}(\mathbf{y}_{k,n}, \mathbf{z}_m)$
- Global score recursion
$$D[m, n] = d[m, n] + \min(D[m-1, n], D[m, n-1], D[m-1, n-1])$$
$$\text{Match}(W_k, S) = -D[M, N]$$

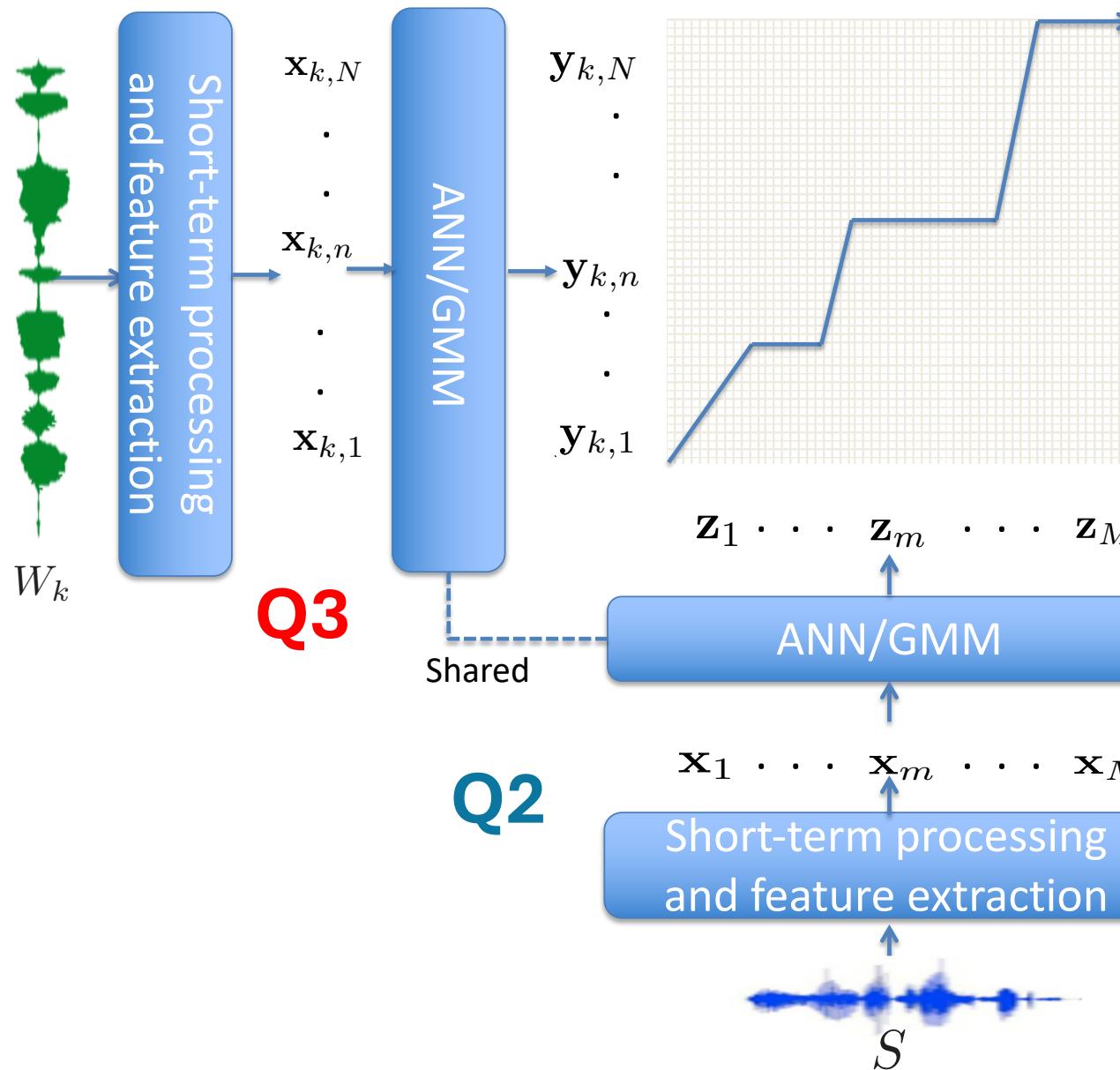
Probabilistic symbols sequence matching



Posterior feature-based speech processing (1)



Posterior feature-based speech processing (2)



Representation of word hypothesis by acoustic instance(s)

Each time frame can be regarded as an HMM state

Text-based representation and acoustic instance-based representation of word hypothesis are interchangeable

Q2: \mathbf{z}_m probability estimator training

Viterbi Expectation-Maximization (EM)

$$Z = \mathbf{z}_1 \cdots \mathbf{z}_m \cdots \mathbf{z}_M$$

E-step
Viterbi algo.

Alignment (m, n)

M-step
Re-train probability estimator

S

Trained/
Initialized
Probability
estimator

Repeat until
convergence

$$Y_k = \mathbf{y}_{k,1} \cdots \mathbf{y}_{k,n} \cdots \mathbf{y}_{k,N}$$

based on W_k

$$\mathbf{y}_{k,n}$$

$$\text{Error} = \text{KL}(\mathbf{y}_{k,n}, \mathbf{z}_m)$$

Gradient

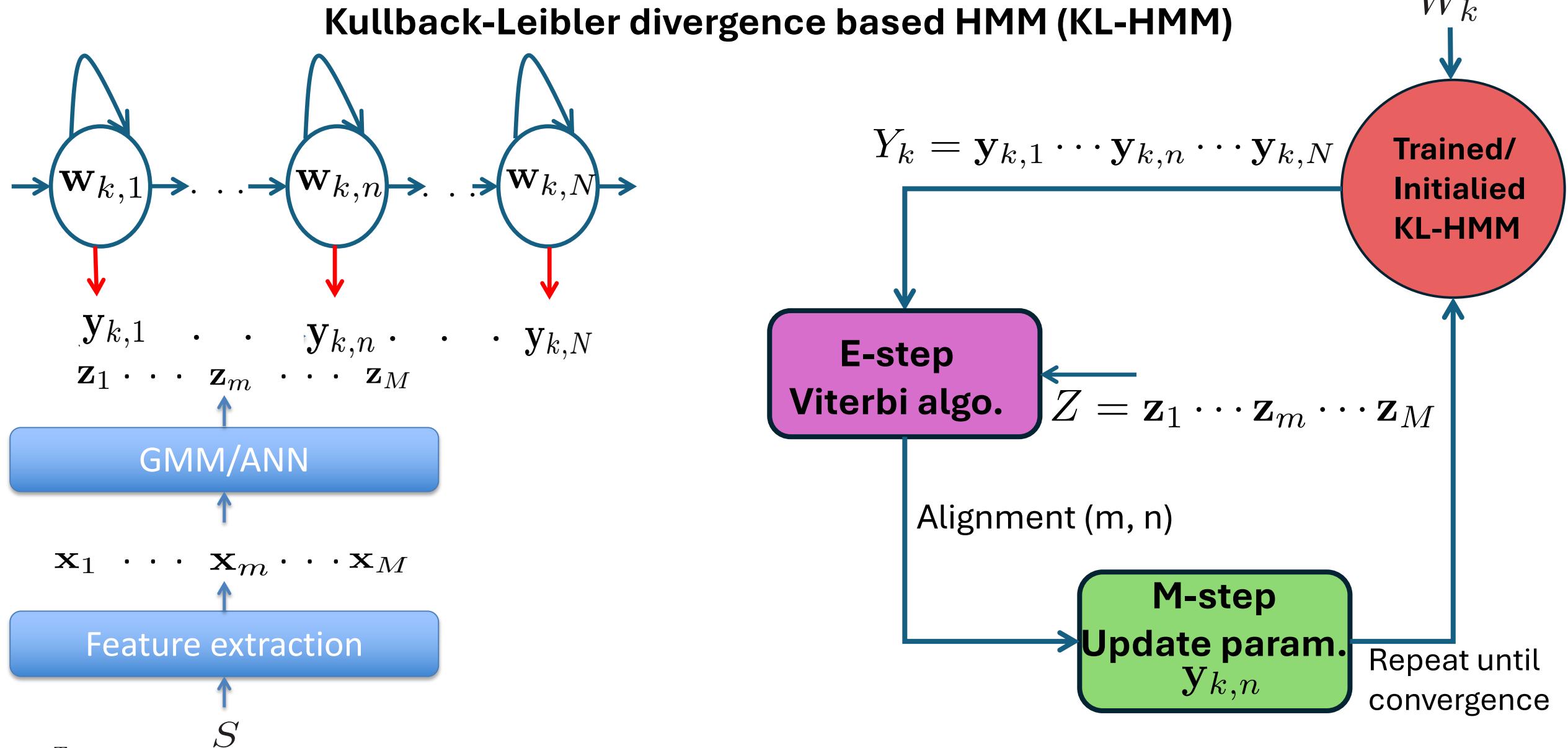
$$\mathbf{z}_m$$

ANN

$$\mathbf{s}_m$$

If $\mathbf{y}_{k,n} = \delta^d$ (Kronecker delta)
 $\text{KL}(\mathbf{y}_{k,n}, \mathbf{z}_m) = -\log(z_m^d) = \text{Cross entropy}$

Q3: $y_{k,n}$ parameters estimation



Comparison of approaches

	Q1	Q2	Q3	Q4	Match()
HMM-based Text only W_k	Phones or Graphemes	Likelihood estimation GMM or ANN	Deterministic	Dyn. Prog.	$p(W_k, S)$
End-to-end Text only W_k	Phones or Graphemes	Probability estimation ANN	Deterministic	Dyn. Prog.	$P(W_k S)$
Posterior feature-based	Different possibilities	Probability estimation Many possible estimators	Probabilistic	Dyn. Prog.	$\frac{p(S W_k)}{p(S \neg W_k)}$
Instance-based (spec.)	Spectral feat. based	Spec. feat. vector extract.	Spec. feat. vector extract.	Dyn. Prog.	$\frac{p(S W_k)}{p(S \neg W_k)}$

Posterior feature-based ASR Demo



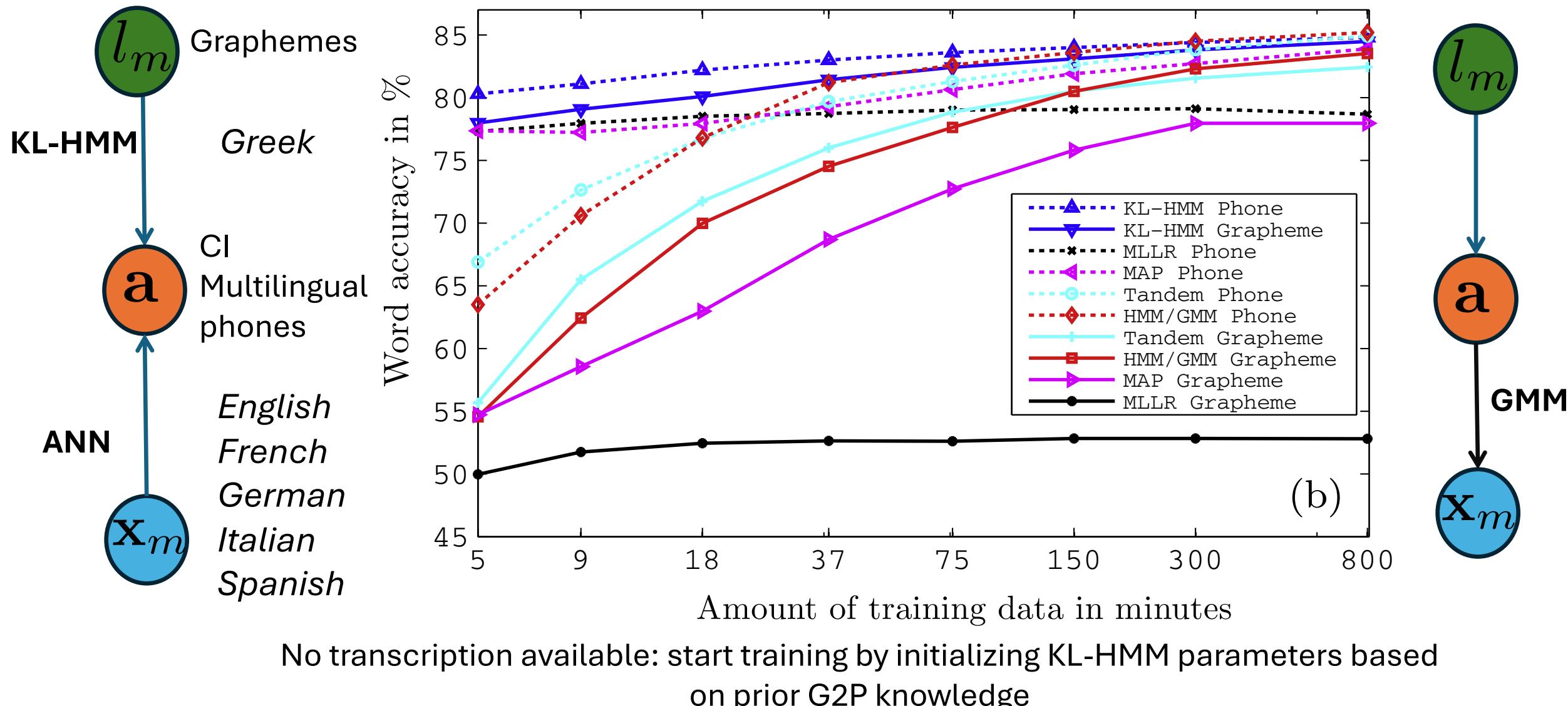
- Instance-based approach
- Posterior feature extractor
 - Two hidden layer MLP trained to classify context-independent phones
- Input: nine frames of 39 dimensional cepstral feature vector (13 static + 13 delta + 13 delta-delta)
- Cross entropy loss
- Stochastic gradient descent

[Posterior Features for Template-based ASR](#), [Serena Soldo](#), [Mathew Magimai-Doss](#), [Joel Praveen Pinto](#) and [Hervé Bourlard](#),

in: Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing, Prague, 2011

[Synthetic References for Template-based ASR using Posterior Features](#), [Serena Soldo](#), [Mathew Magimai-Doss](#) and [Hervé Bourlard](#), in: Proceedings of Interspeech, 2012

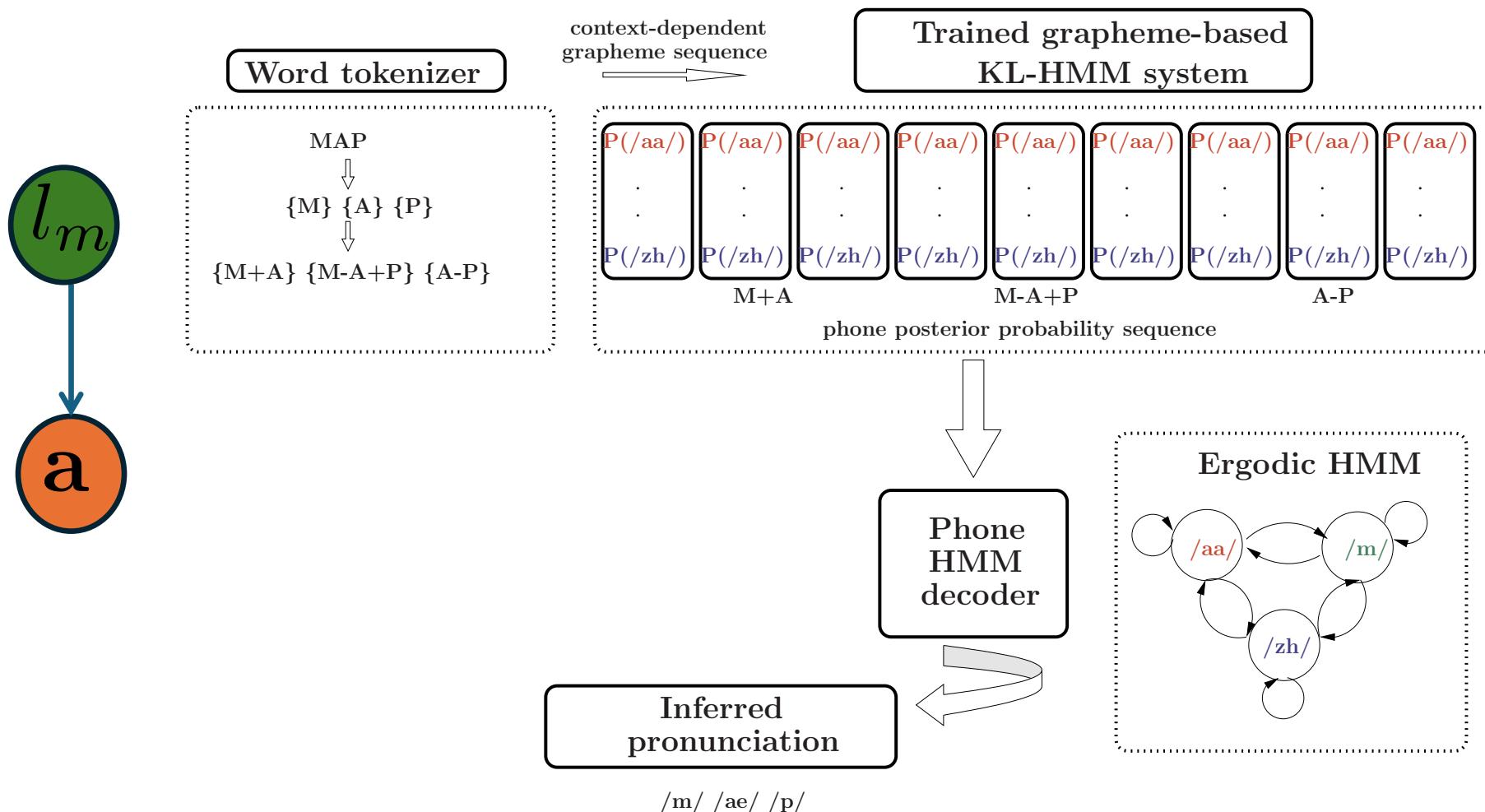
Handling acoustic and lexical resource constraints



[Acoustic and Lexical Resource Constrained ASR using Language-Independent Acoustic Model and Language-Dependent Probabilistic Lexical Model](#), Ramya Rasipuram and Mathew Magimai.-Doss, in: Speech Communication, 68:23–40, 2015

[Probabilistic Lexical Modeling and Unsupervised Training for Zero-Resourced ASR](#), Ramya Rasipuram, Marzieh Razavi and Mathew Magimai.-Doss, in: Proc. of ASRU, 2013

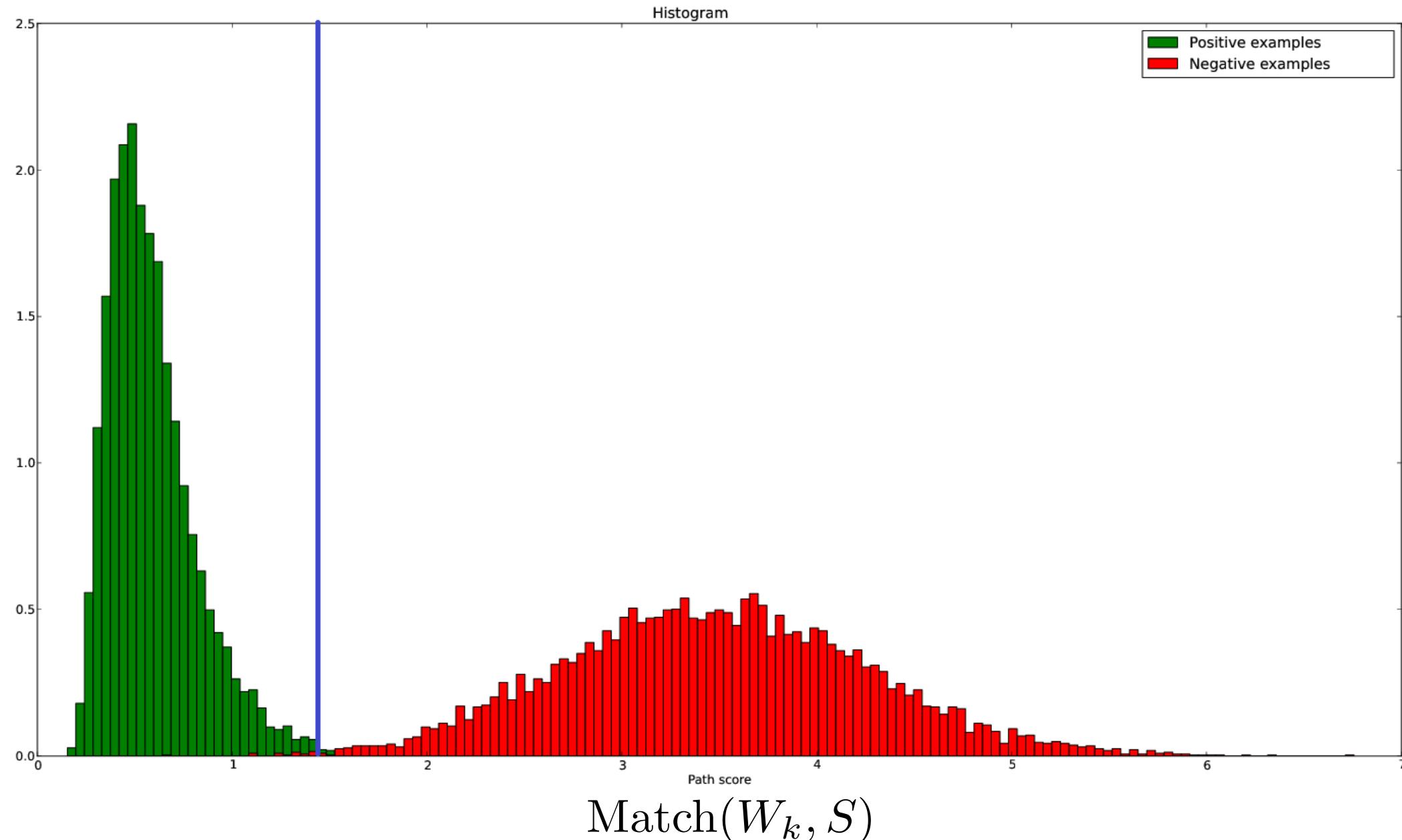
Lexical resource development



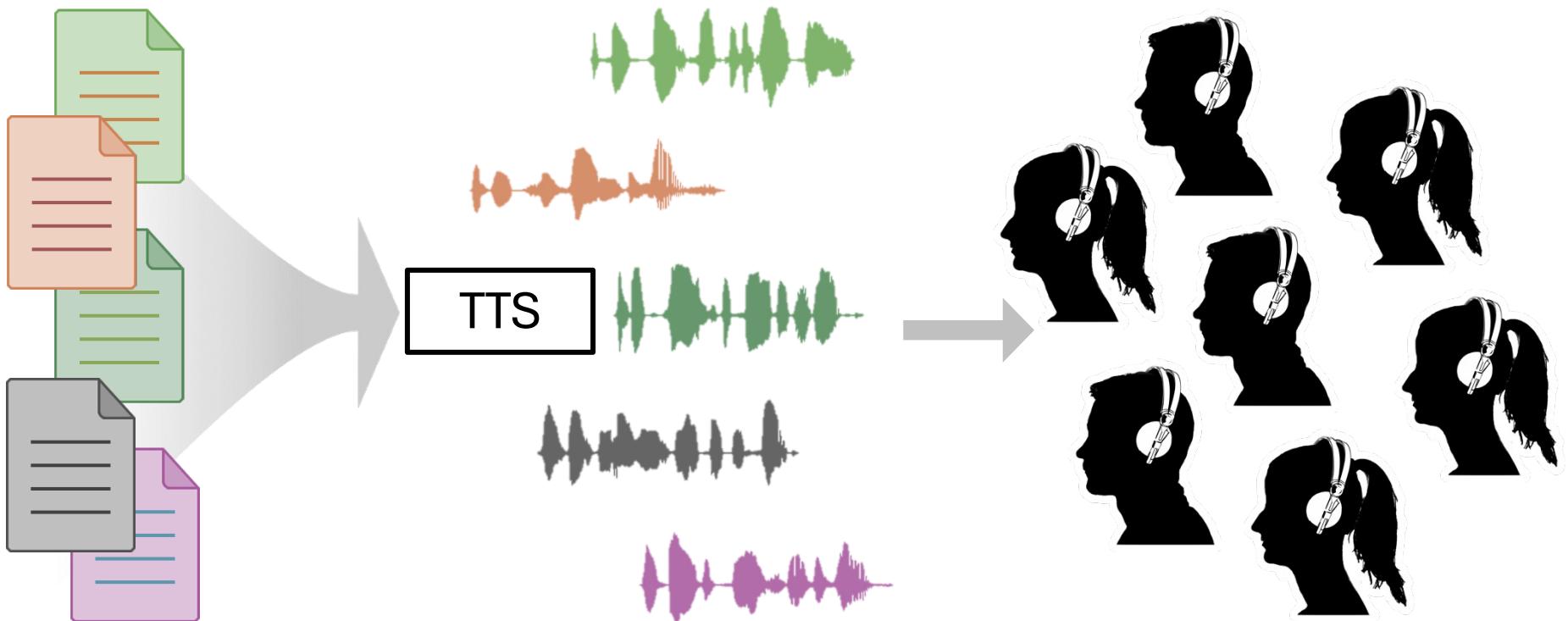
- Can be extended to automatic subword derivation and lexicon development
- Combination of multiple G2P converters
- Unified framework for acoustic-to-phone conversion and G2P conversion

[Acoustic data-driven grapheme-to-phoneme conversion in the probabilistic lexical modeling framework](#), M. Razavi, R. Rasipuram and M. Magimai.-Doss, in: Speech Communication, 80, 2016
[Towards Weakly Supervised Acoustic Subword Unit Discovery and Lexicon Development Using Hidden Markov Models](#), M. Razavi, R. Rasipuram and M. Magimai-Doss, in: Speech Communication, 96:168-183, 2018
M. Razavi, [On Modeling the Synergy Between Acoustic and Lexical Information for Pronunciation Lexicon Development](#), PhD Thesis 7851, EPFL, 2017

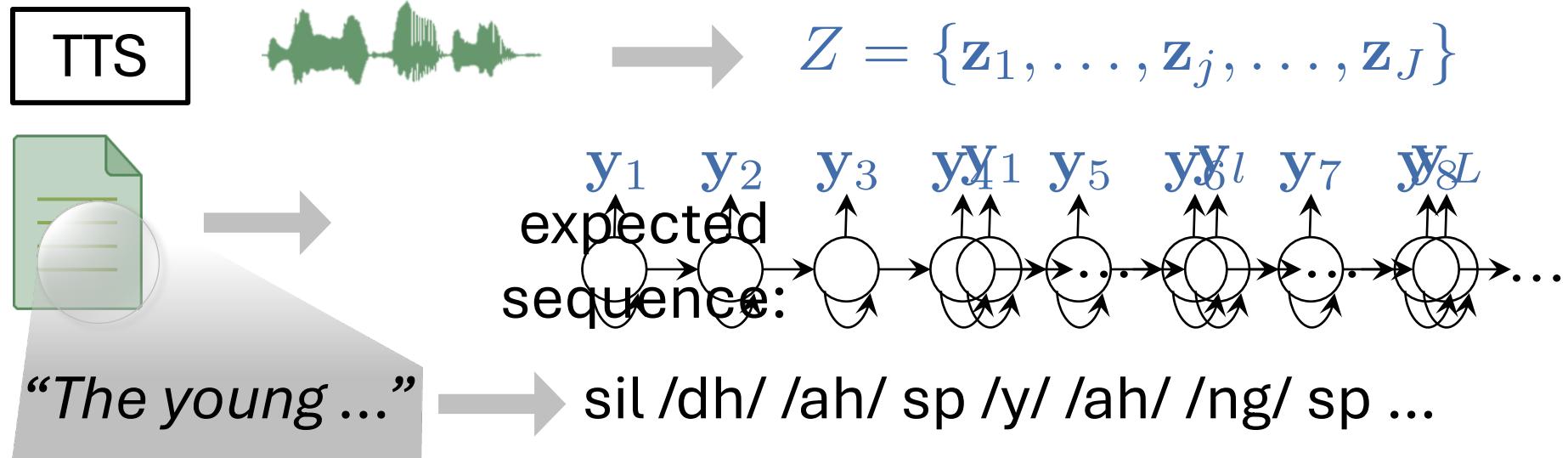
Probabilistic symbols sequence matching



Subjective intelligibility assessment (1)

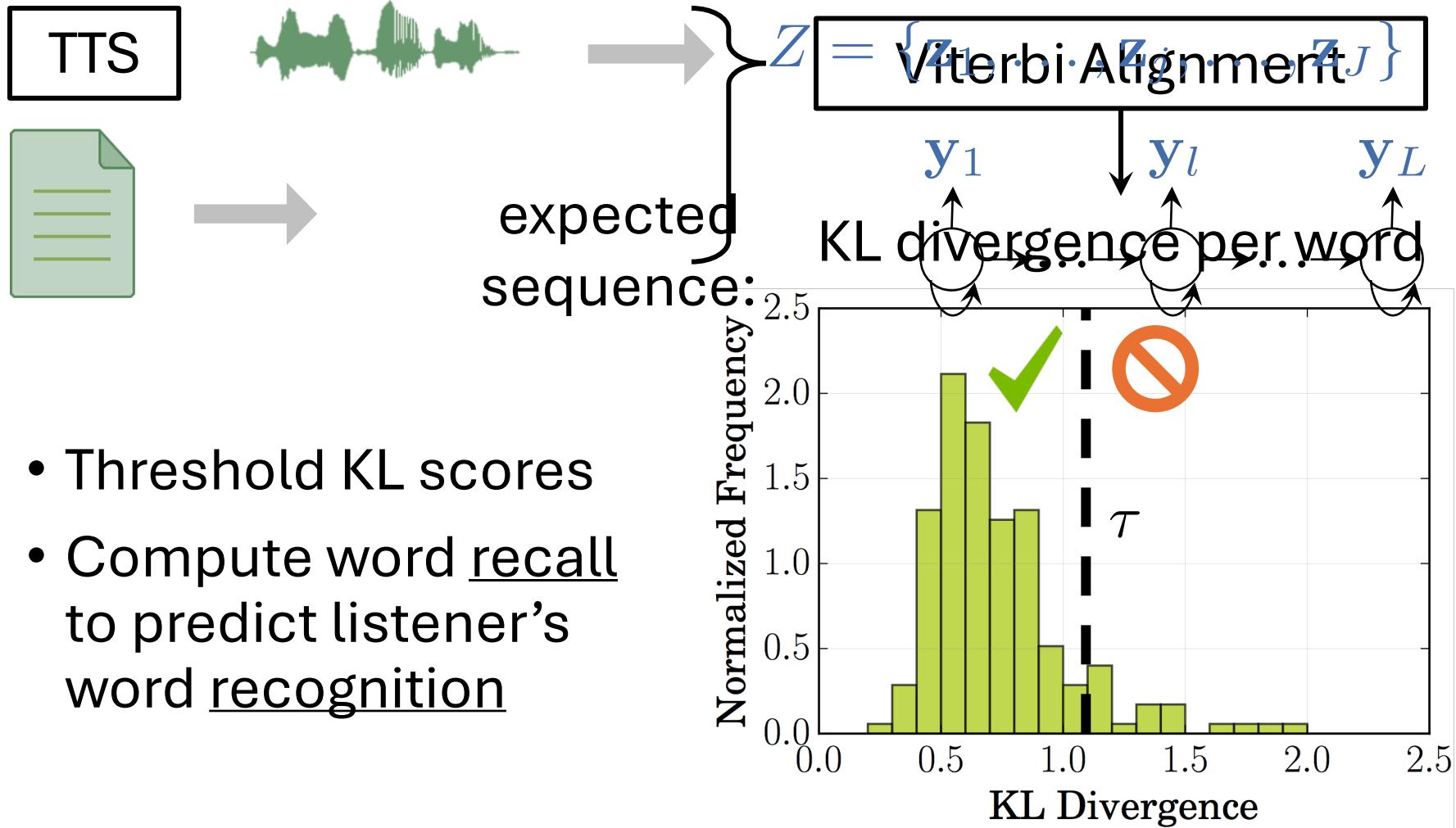


Match word hypothesis with synthetic speech



- The phonetic transcription is obtained from a lexicon
- HMM states are parameterized by distributions y_l , with the KL divergence $\text{KL}(y_l, z_j)$ as local score

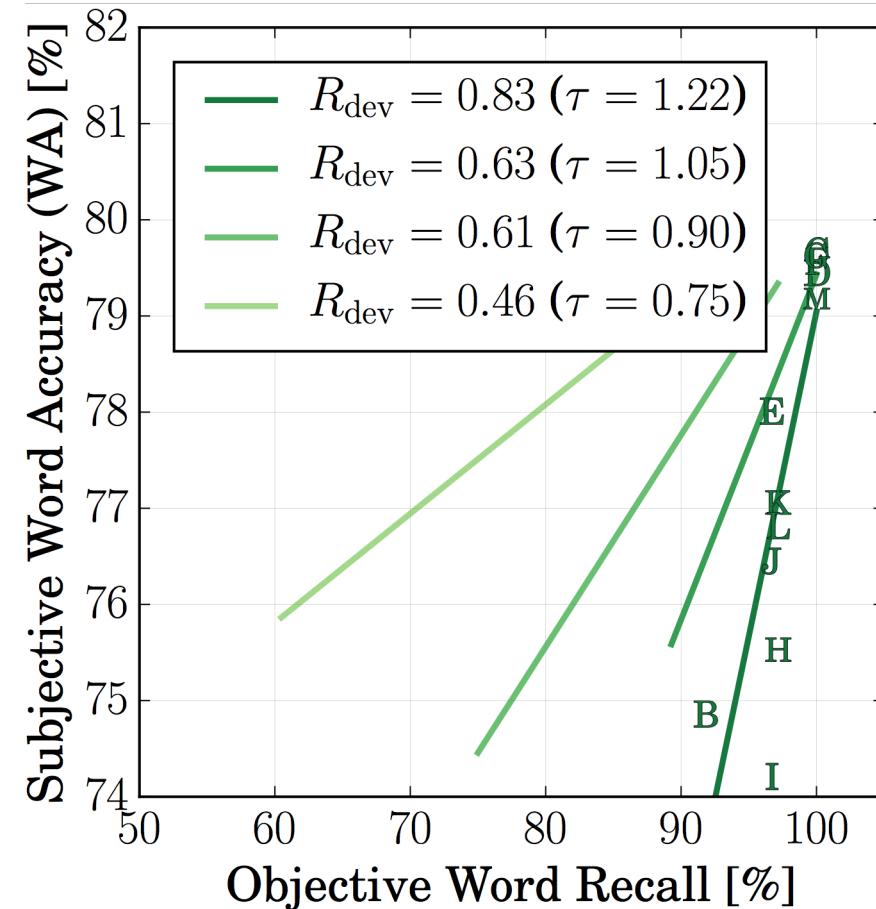
Hypothesis testing and word recall



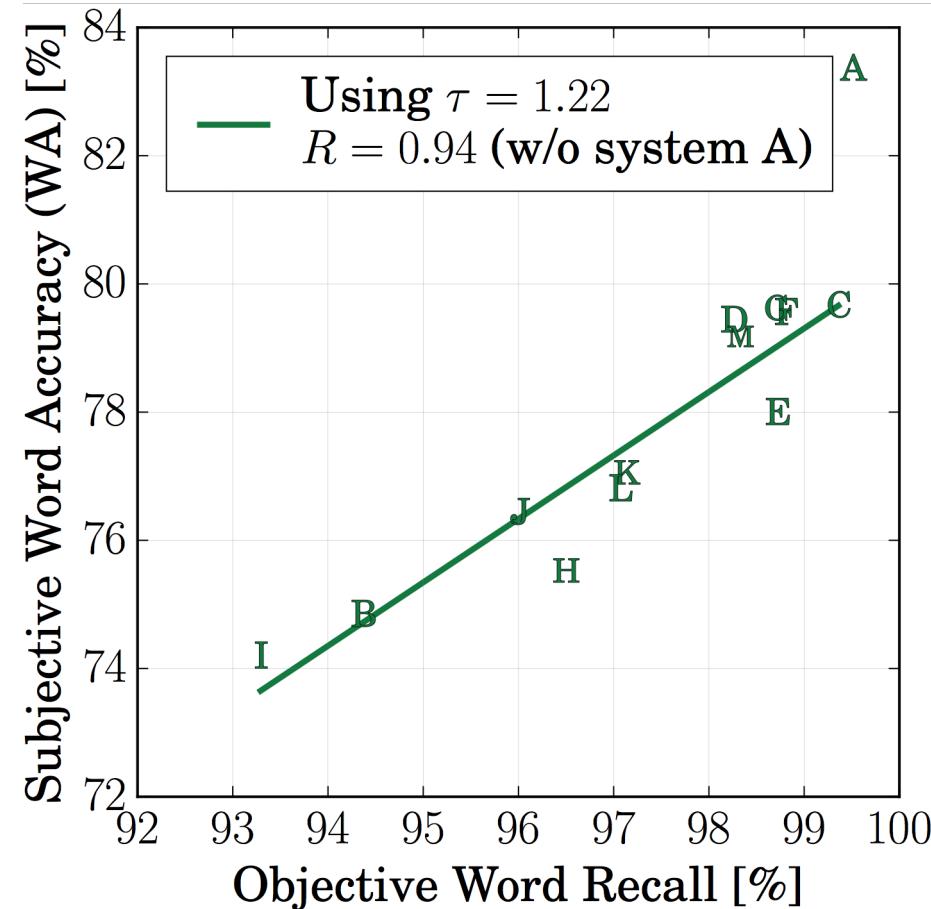
Evaluation

2011 Blizzard Challenge

- 12 TTS systems named “B” to “M”
- 26 semantically unpredictable sentences (SUS), scored by 231 listeners (26 scores per listener)
- “System A”: natural speech from a voice talent

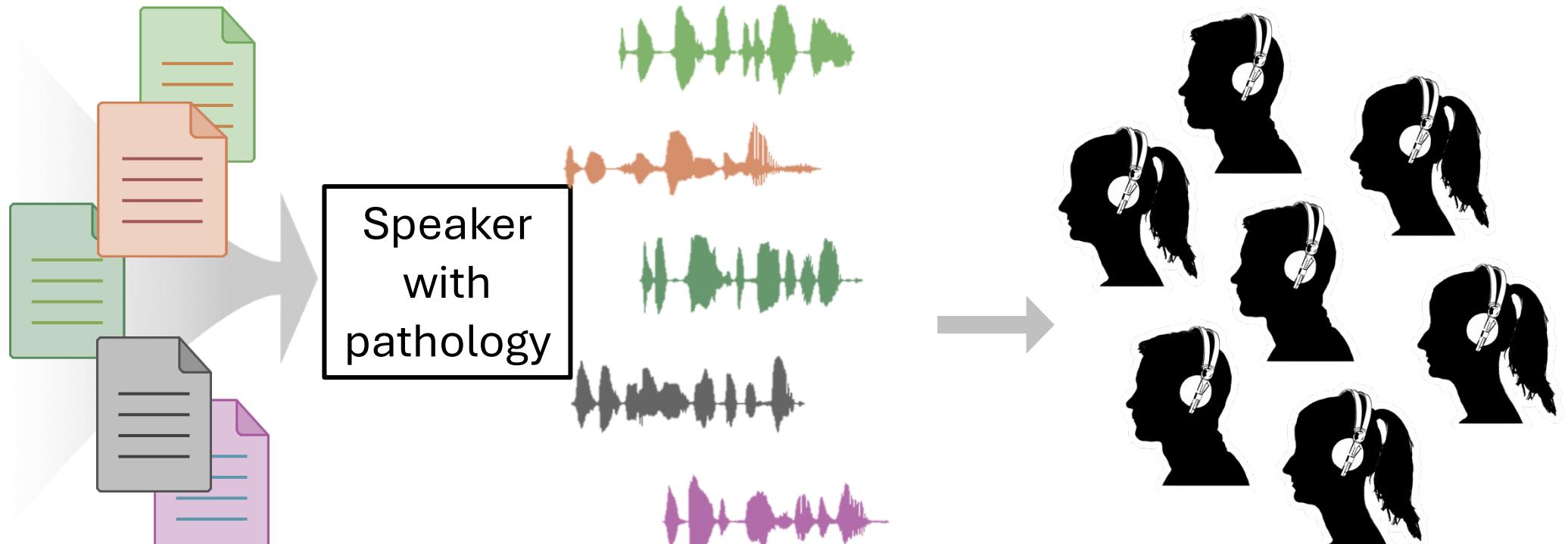


Determine threshold on development data



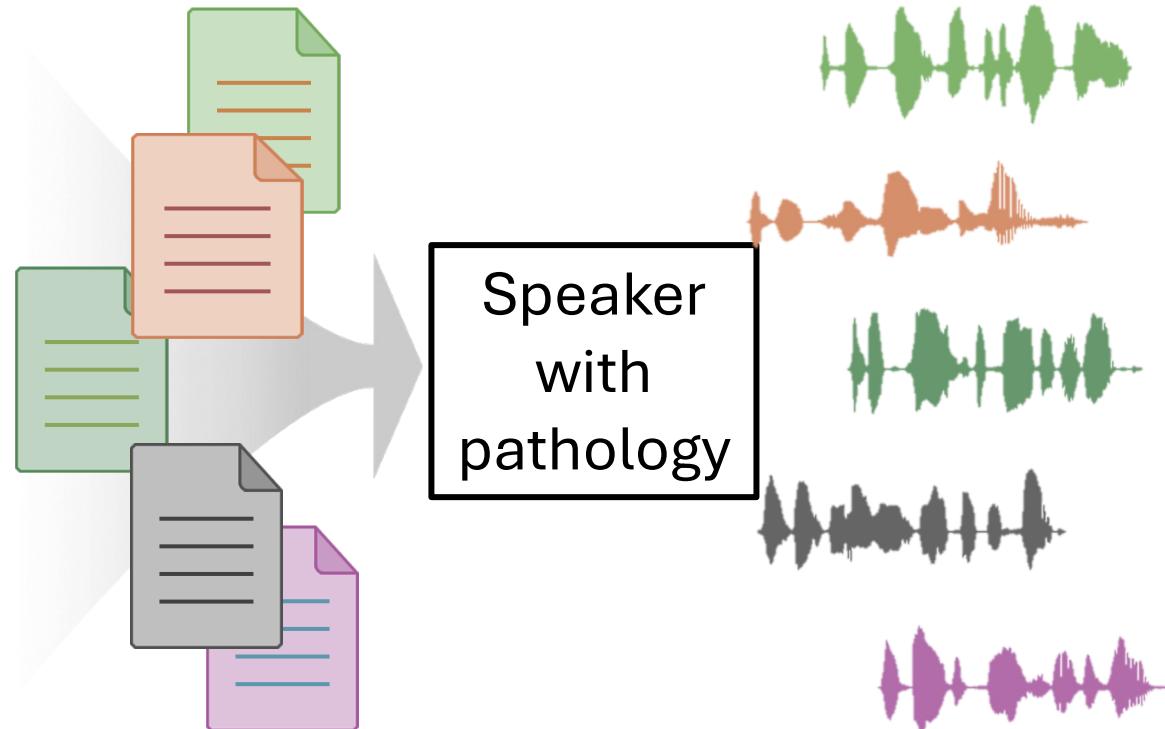
Performance on evaluation data

Subjective intelligibility assessment (2)

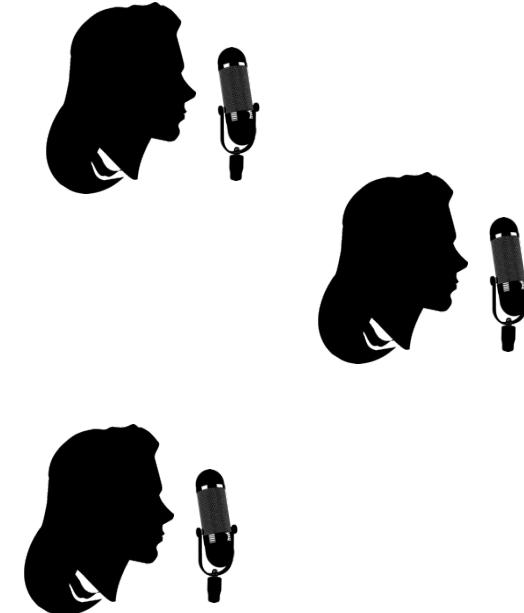


e.g., speaker with dysarthria

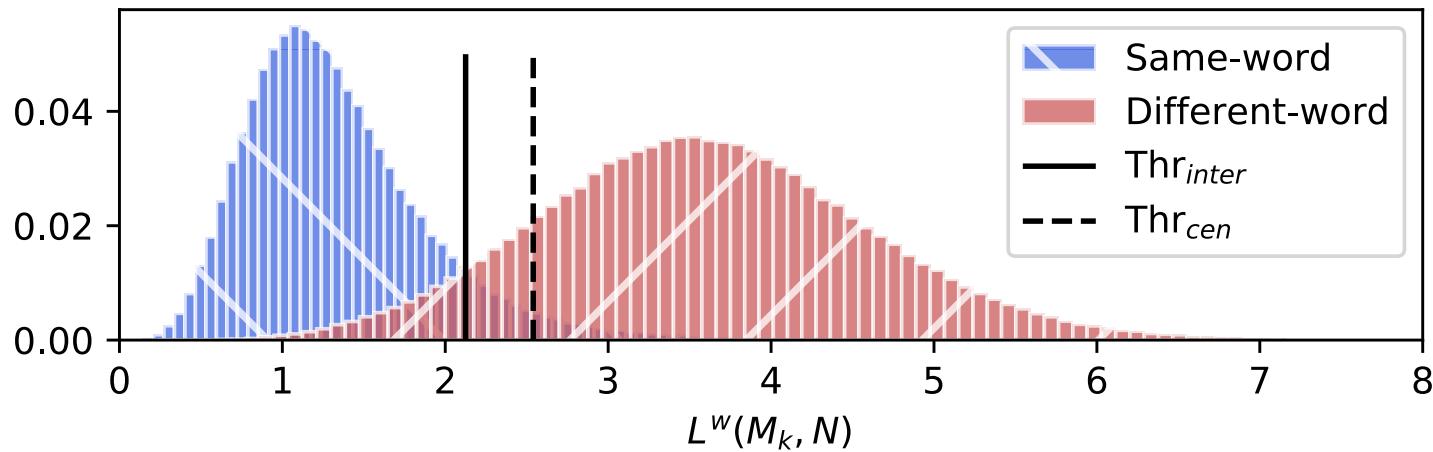
Objective intelligibility assessment (1)



Replace listeners by healthy control speakers' speech utterances of each word or even synthetic speech.



Objective intelligibility assessment (2)

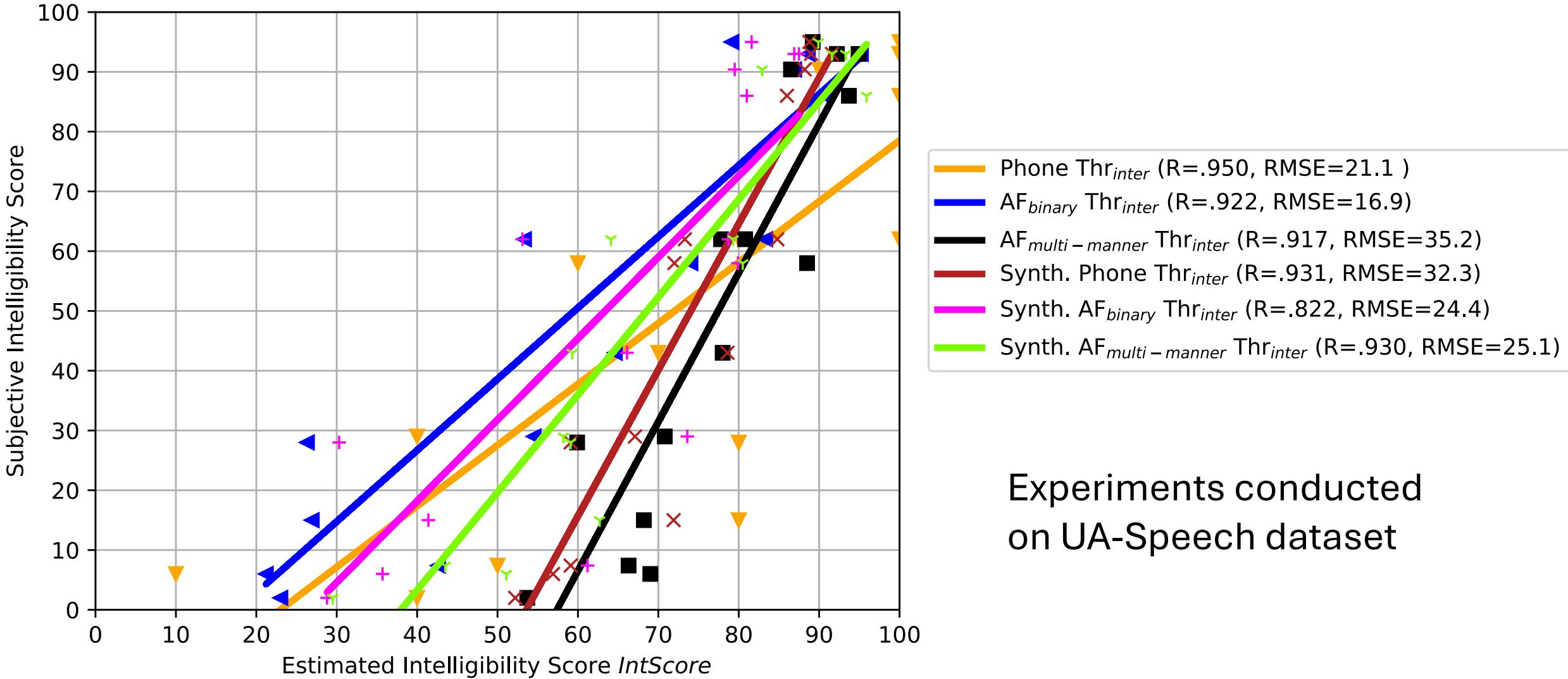


Determining threshold from healthy control speakers' data by forming (a) same-word pair and (b) different word pair.

For each pathological speech utterance

1. Performance utterance by matching to the respective word utterances from healthy control speakers
2. Count number of same word decisions
3. Count more than or equal to half of number of healthy control speakers then decide the pathological speech utterance is recognizable else non-recognizable
4. Count the number of recognizable words

Objective intelligibility assessment (3)



Posterior feature based approach

A modular and flexible framework that enables abstraction and integrated modeling of communication phenomenon

- **Q1** – Flexibility to choose and/or explore symbol set
- **Q2** – Flexibility to choose symbol probability estimation method
- **Q3** – Flexibility to text-based and acoustic instance-based representation of word hypothesis
- **Q4** – Interpretable and explainable information-theoretic approach akin to string matching

Systematically assimilate developments in different fields such as, linguistics, signal processing, machine learning, information theory, coding theory, communication theory, and probability and statistics

