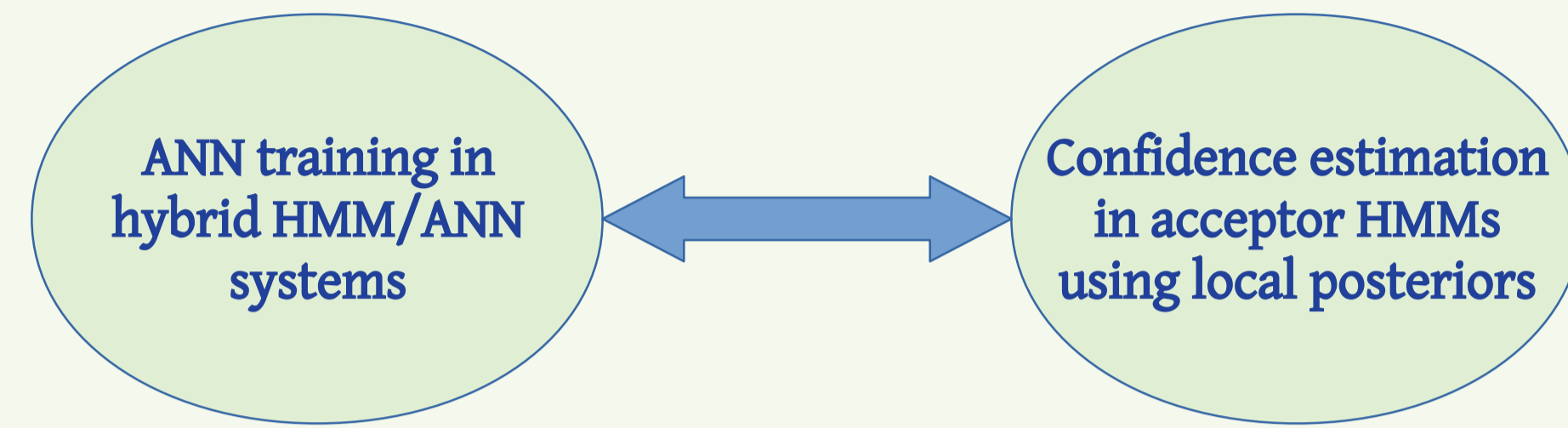
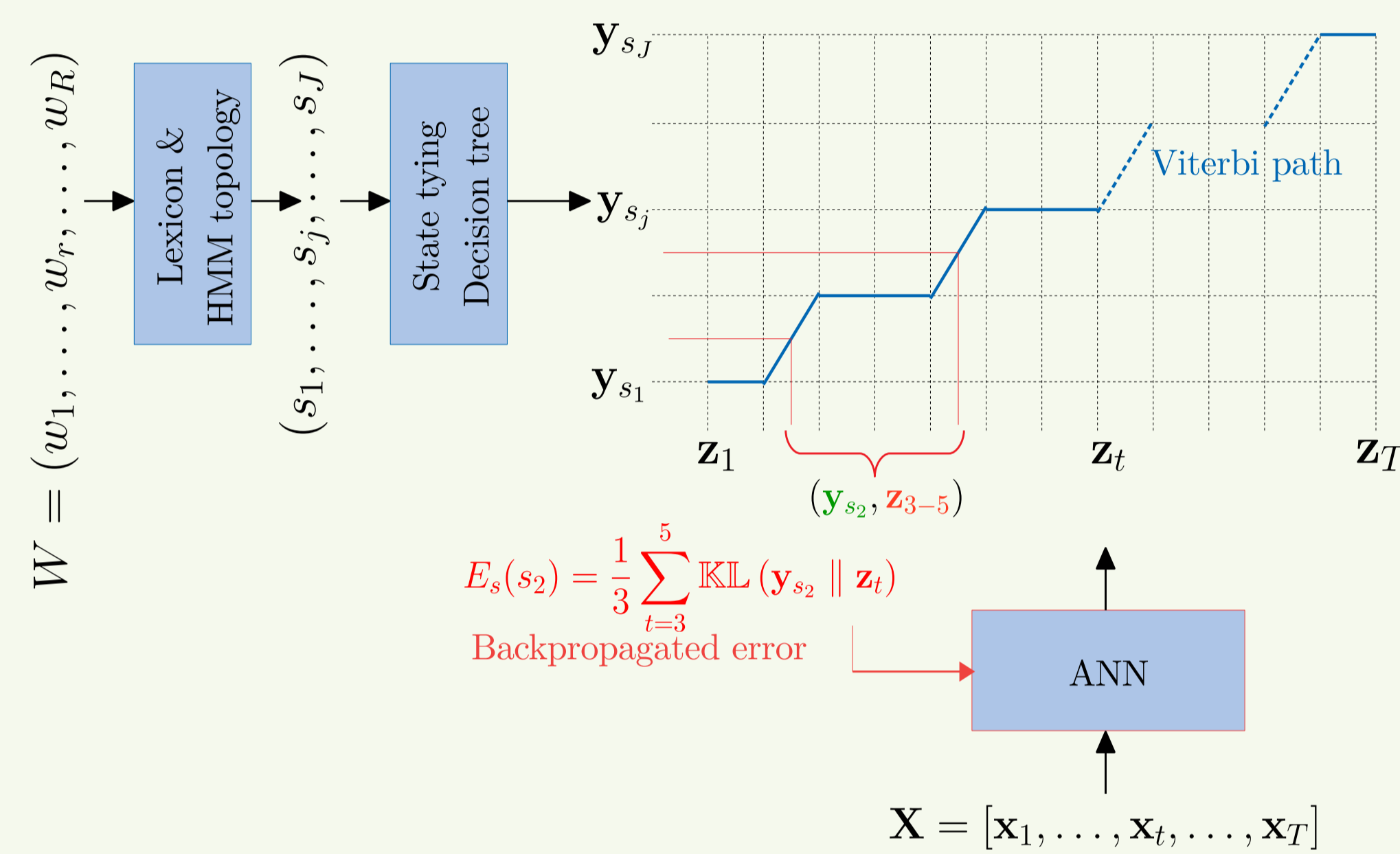


Segment-level training based on Confidence Measures for Hybrid HMM/ANN Speech Recognition

1 In this Work



4 Proposed Segment-level ANN Training



- When the segments represent HMM states, a state-level error $E_s(s_j)$ can be minimised:

$$E_s(s_j) = -CM(s_j) = \frac{\sum_{t=b(s_j)}^{e(s_j)} \mathbb{KL}(y_{s_j} \parallel z_t)}{e(s_j) - b(s_j) + 1}$$

- When the segments represent phone units, a phone-level error $E_{ph}(ph_k)$ can be minimised:

$$E_{ph}(ph_k) = \frac{1}{N_{ph_k}} \sum_{n=1}^{N_{ph_k}} E(s_{j+n}); \text{ where } ph_k \text{ comprises } N_{ph_k} \text{ states}$$

- Priors $P(a^d)$ are estimated from the state segment counts instead of frame label counts.

5 Data Sets and Experimental Setup

	AMI	Mediaparl German	Mediaparl French	TIMIT
Training hours	77.3	14.5	16.1	3.1
Phone set count	176	57	38	48
Vocabulary size	52.5k	16.7k	12.4k	48
LM order	3-gram	2-gram	2-gram	2-gram
Features	fMLLR+spk-iVec	MFCC		
Tools	Kaldi+Keras/Tensorflow			
Alignments	From HMM-SGMM systems			
Training	E_f, E_s or E_{ph} , followed by sMBR			

2 Background: Conventional Hybrid HMM/ANN Systems

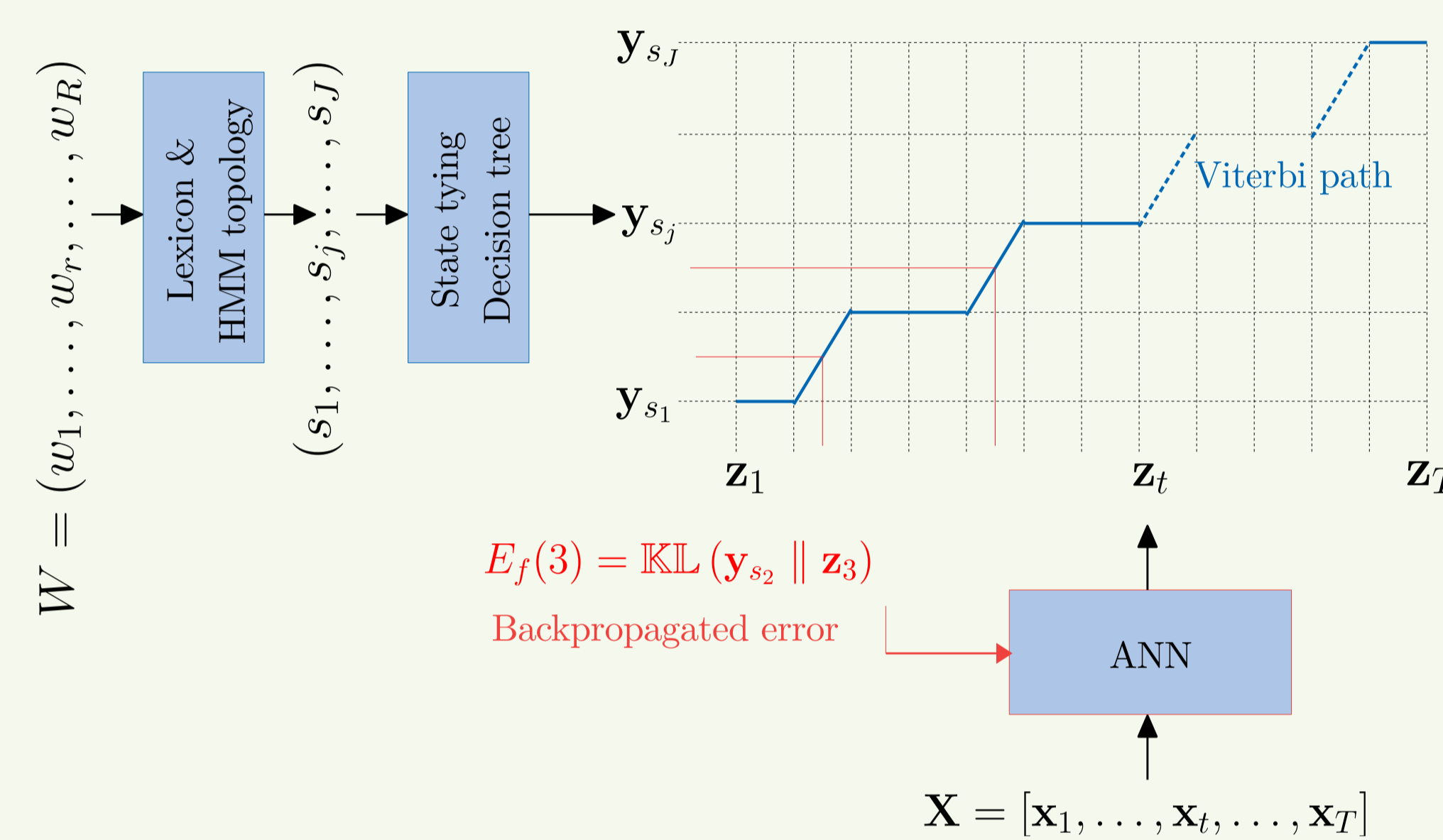
- In hybrid hidden Markov model (HMM) based speech recognition, the *scaled likelihood* of an acoustic observation x_t given a HMM state q_t at time t , labelled l^i , is estimated as:

$$\frac{p(x_t | q_t = l^i)}{p(x_t)} = \sum_{d=1}^D \frac{p(x_t | a^d)}{p(x_t)} P(a^d | q_t = l^i) = \sum_{d=1}^D \frac{P(a^d | x_t)}{P(a^d)} P(a^d | q_t = l^i)$$

↑ priors ↑ decision tree

- Conventionally, given the *segmentation*, the artificial neural network (ANN) is trained using one hot encodings of the targets and minimising frame level cross-entropy. This can be expressed as:

$$E_f(t) = \mathbb{KL}(y_{s_j} \parallel z_t) \quad (s_j=l^j) \rightarrow a^{d'} \quad \mathbb{KL}(\delta_{d'} \parallel z_t) = -\log(P(a^{d'} | x_t)) \quad (1)$$



6 Results

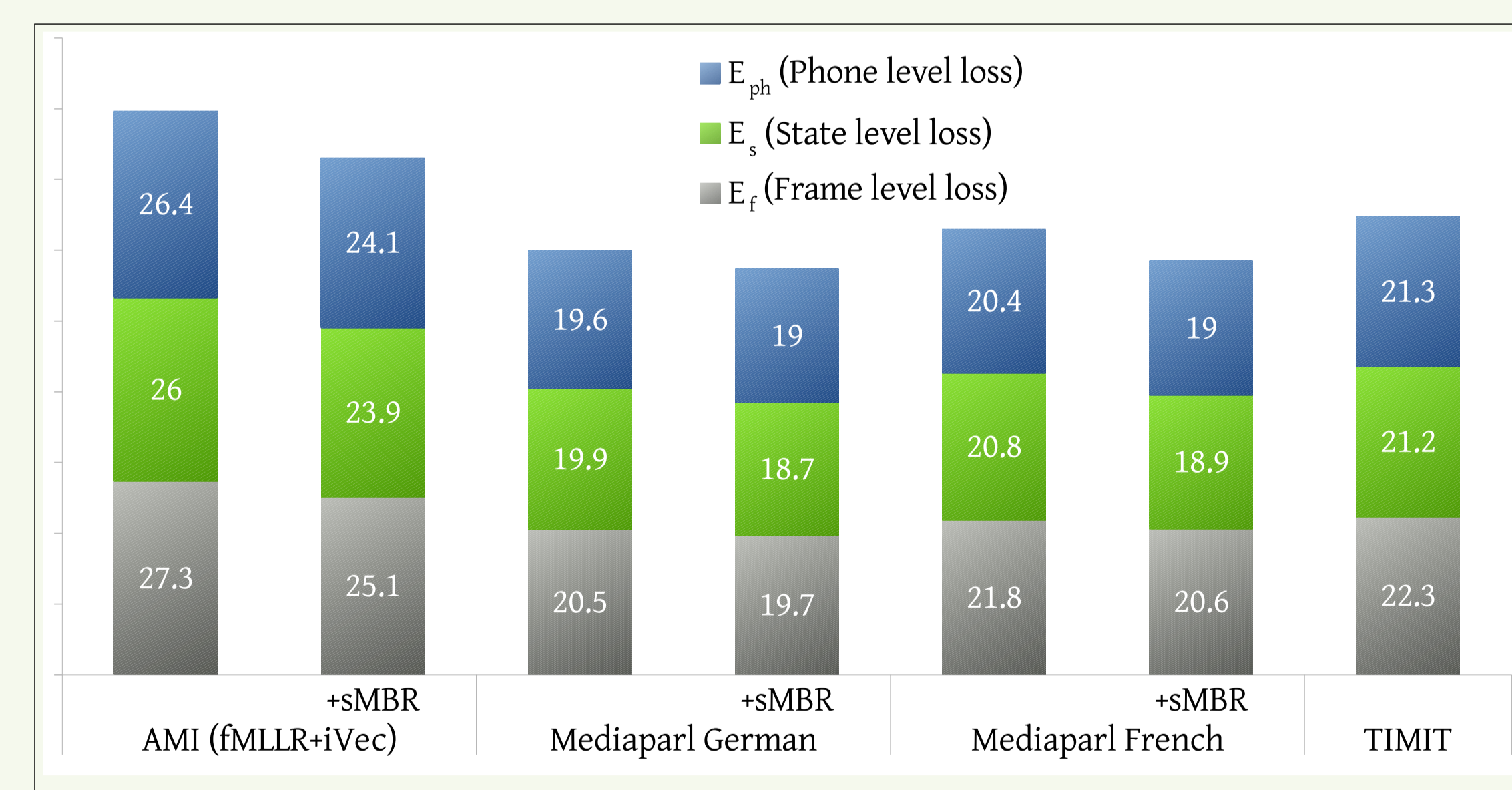


Figure: Word error rate on AMI and Mediaparl data sets, phone error rate on TIMIT

Acknowledgements

HASLERSTIFTUNG

This work was funded by the Hasler foundation under the project Flexible Linguistically-guided Objective Speech aSessment (FLOSS) and was entirely carried out at the Idiap Research Institute.

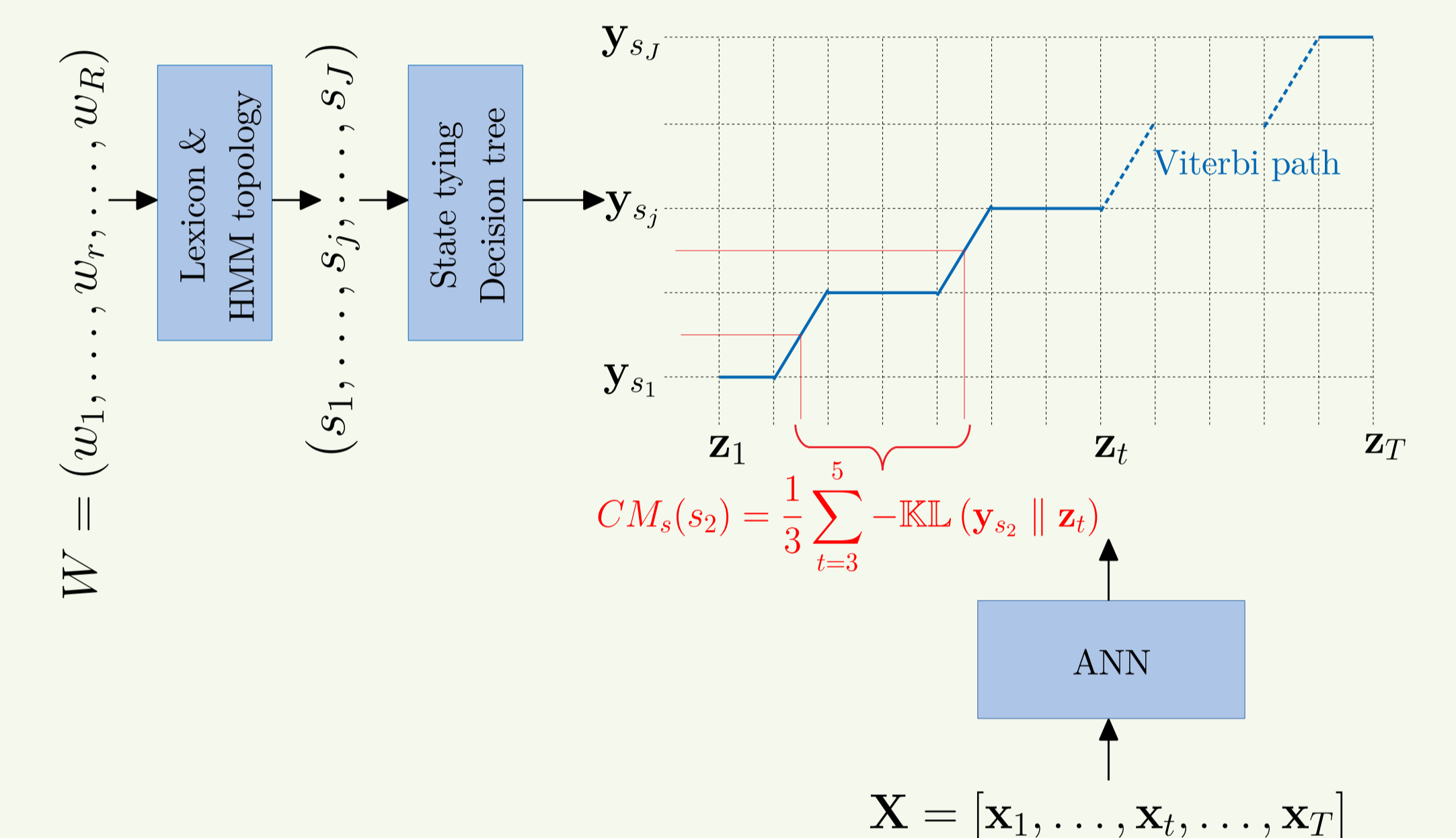
3 Confidence Estimation using Local Posteriors

- Given an alignment between X and W and the local posterior probability estimates, a confidence measure $CM(s_j)$ can be estimated by rescaling the segment of s_j at $t = \{b(s_j), \dots, e(s_j)\}$ as

$$CM(s_j) = \frac{\sum_{t=b(s_j)}^{e(s_j)} \log(P(q_t = l^j | x_t))}{e(s_j) - b(s_j) + 1} \stackrel{l^j \rightarrow a^{d'}}{=} \frac{\sum_{t=b(s_j)}^{e(s_j)} \log(P(a^{d'} | x_t))}{e(s_j) - b(s_j) + 1}$$

- Based on (1), we can express the state level confidence $CM(s_j)$ estimation as a matching of Y and Z with a local cost based on Kullback-Leibler (KL) divergence:

$$CM(s_j) = \frac{\sum_{t=b(s_j)}^{e(s_j)} -\mathbb{KL}(y_{s_j} \parallel z_t)}{e(s_j) - b(s_j) + 1}$$



7 Analysis: Effect of Silence Duration on the Training

- Here we show how increasing the silence duration affects training using the three cost functions.
- Silence was artificially added at the beginning and end of each training and test utterance.

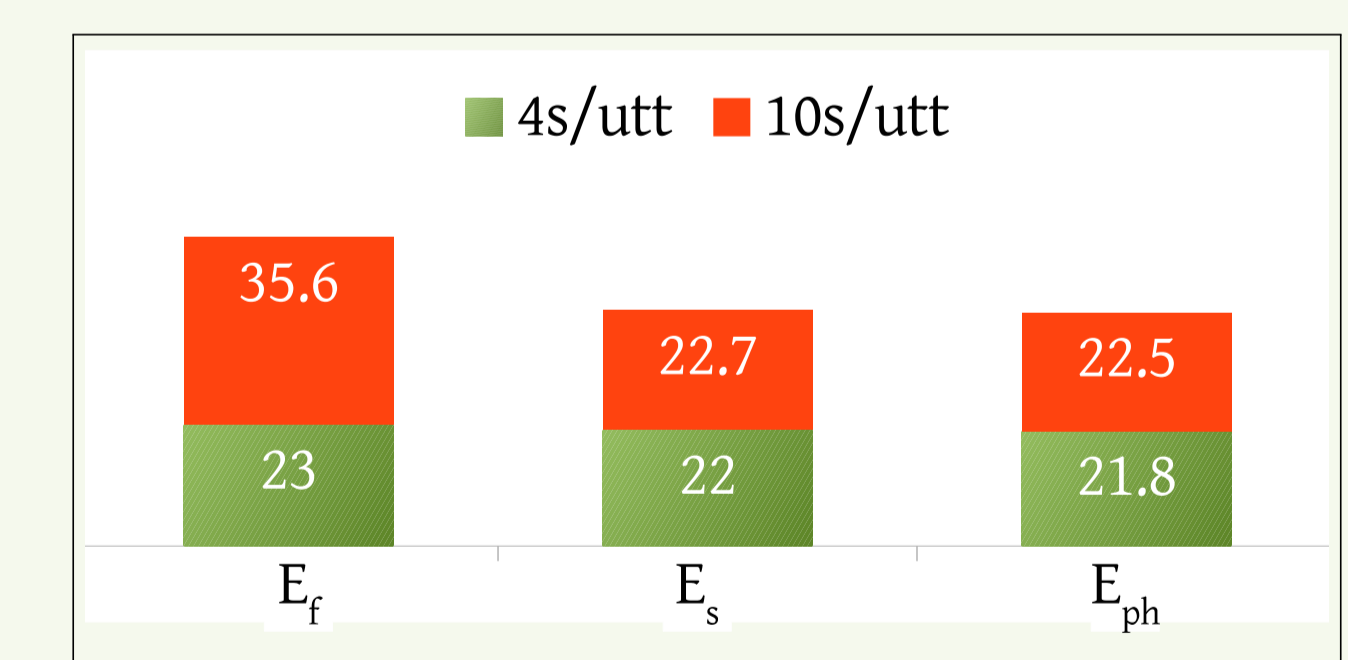


Figure: Phone error rate on TIMIT

8 Conclusion

The proposed linguistic-segment-level training of ANNs based on confidence measures

- yields better systems than using frame-level cross-entropy.
- adds to the efficacy of further sequence discriminative training.
- improves robustness to duration variations in the training data set.

