# Confidence Measures for CTC-based Phone Synchronous Decoding

**Zhehuai Chen**, Yimeng Zhuang, Kai Yu

- **ASR Decoding**
  - Inference by AM/LM/lex …
  - Model and search are both imperfect



Speech Input

$S$

Speech Analysis

$O$

Pronunciation Lexicon

Acoustic Model $p(O|W)$

Decoder
$\hat{W} = \underset{w \in W}{\text{argmax}}\ p(O|W)P(W)$

$P(W)$ Language Model

$\hat{W}$

Recognition Result

- **ASR Decoding**
  - Inference by AM/LM/lex …
  - Model and search are both imperfect

- **Confidence Measure (CM)**
  - Reliability evaluation of ASR results
  - Traditional CM
    - Predictor features based CM
      - Acoustic score, duration, entropy … (NOT ideal)
      - CRF, NN … (need training stage; train ≠ test)
    - Hypothesis Posterior based CM
      - Theoritically sounder



Speech Input

$\downarrow$ S

Speech Analysis

$\downarrow$ O

Pronunciation Lexicon

Acoustic Model

$p(O|W)$

Decoder

$\hat{W} = \underset{w \in W}{\arg\max} \; p(O|W)P(W)$

Language Model

$P(W)$

$\downarrow \hat{W}$

Recognition Result

- ASR as the *maximum a posterior* (MAP) decision process

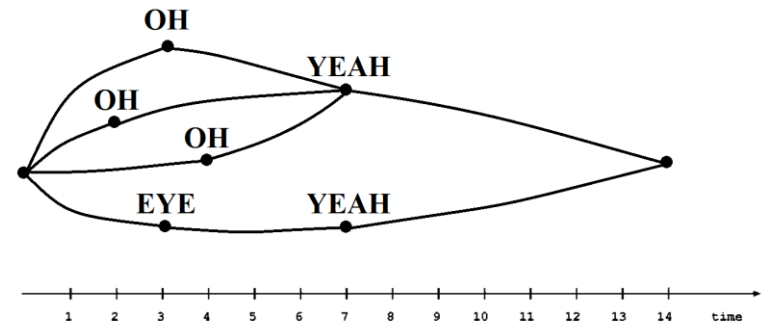$$\widehat{W} = \arg\max_{W \in \Sigma} p(W \mid X)$$

$$= \arg\max_{W \in \Sigma} \frac{p(X \mid W) \cdot p(W)}{p(X)}$$

$$p(X) = \sum_H p(X, H) = \sum_H p(H) \cdot p(X \mid H)$$

- H is from lattice/filler
  - Both imperfect

3

- ASR as the *maximum a posterior* (MAP) decision process

$$\widehat{W} = \arg\max_{W \in \Sigma} p(W \mid X)$$

$$= \arg\max_{W \in \Sigma} \frac{p(X \mid W) \cdot p(W)}{p(X)}$$

$$p(X) = \sum_H p(X, H) = \sum_H p(H) \cdot p(X \mid H)$$

- H is from lattice/filler
  - Both imperfect

- **Lattice quality** is the bottleneck
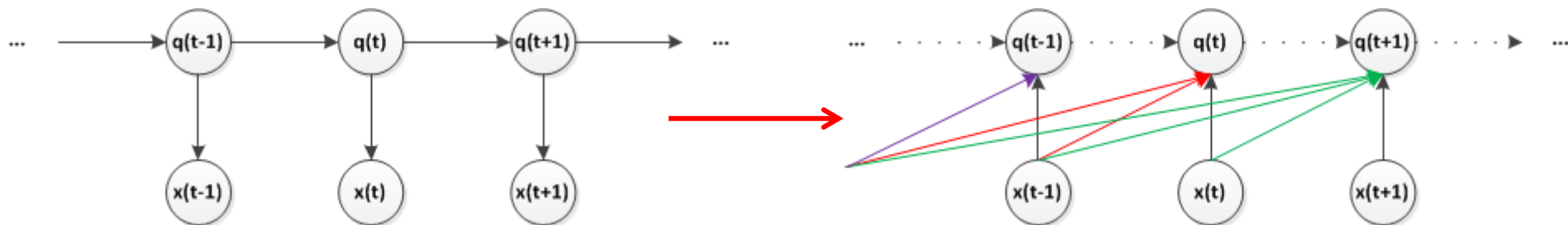  - Not compact
    - Boundary unstable



  - Not precise
    - Beam prune
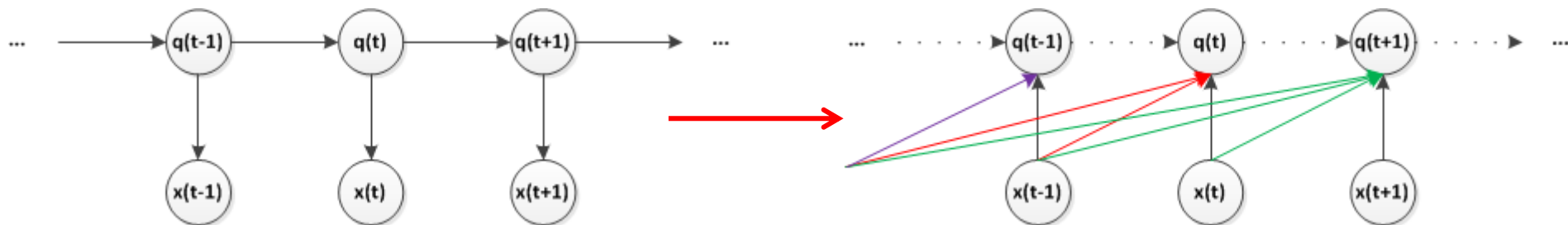
# WE NEED NEW MODEL !

4

■ From HMM to CTC: do better in *sequential modeling*

- From HMM to CTC: do better in *sequential modeling*



- CTC model: learn the many-to-one function of $\mathcal{B}$

$$P(\mathbf{l}|\mathbf{x}) = \sum_{\boldsymbol{\pi} \in \mathcal{B}^{-1}(\mathbf{l})} P(\boldsymbol{\pi}|\mathbf{x}) = \sum_{\pi:\pi \in L', \mathcal{B}(\boldsymbol{\pi}_{1:T})=\mathbf{l}} \prod_{t=1}^{T} y_{\pi_t}^t$$

$$\mathcal{B} : L' \mapsto L$$
$$L' = L \cup \{\texttt{blank}\}$$

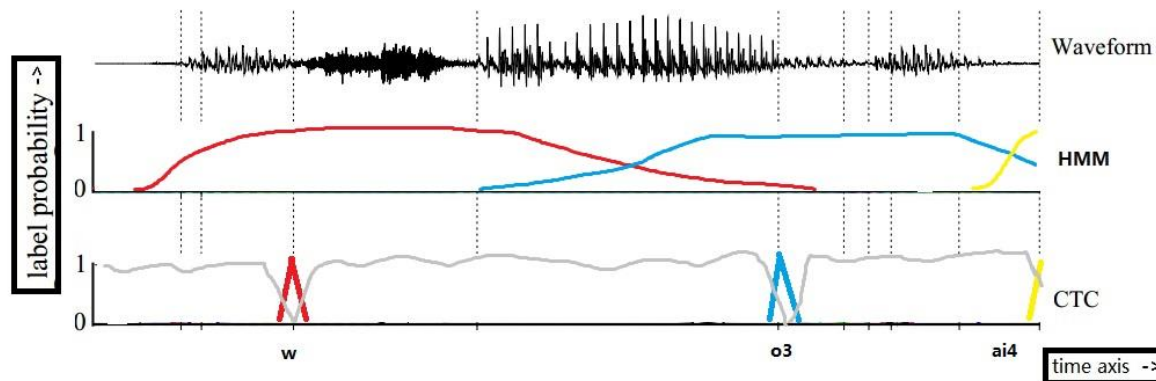- From HMM to CTC: do better in *sequential modeling*



- CTC model: learn the many-to-one function of $\mathcal{B}$

$$P(\mathbf{l}|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{l})} P(\pi|\mathbf{x}) = \sum_{\pi:\pi \in L', \mathcal{B}(\pi_{1:T})=\mathbf{l}} \prod_{t=1}^{T} y_{\pi_t}^t \qquad \mathcal{B} : L' \mapsto L$$
$$L' = L \cup \{\texttt{blank}\}$$

- peaky distribution and concentrated information output

- **frame synchronous Viterbi beam search in CTC**

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}}\{P(\mathbf{w})p(\mathbf{x}|\mathbf{w})\} = \underset{\mathbf{w}}{\operatorname{argmax}}\{P(\mathbf{w})p(\mathbf{x}|\mathbf{l_w})\} \qquad (1)$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}}\left\{P(\mathbf{w})\max_{\mathbf{l_w}}\frac{P(\mathbf{l_w}|\mathbf{x})}{P(\mathbf{l_w})}\right\} \qquad (2)$$

$$\cong \underset{\mathbf{w}}{\operatorname{argmax}}\left\{P(\mathbf{w})\max_{\pi:\pi\in L',\mathcal{B}(\pi_{1:T})=\mathbf{l_w}}\frac{1}{P(\mathbf{l_w})}\prod_{t=1}^{T}y_{\pi_t}^t\right\} \qquad (3)$$

# Frame Sync. to Phone Sync.

- **frame synchronous Viterbi beam search in CTC**

$$\mathbf{w}^* = \operatorname*{argmax}_{\mathbf{w}}\{P(\mathbf{w})p(\mathbf{x}|\mathbf{w})\} = \operatorname*{argmax}_{\mathbf{w}}\{P(\mathbf{w})p(\mathbf{x}|\mathbf{l_w})\} \qquad (1)$$

$$= \operatorname*{argmax}_{\mathbf{w}}\left\{P(\mathbf{w})\max_{\mathbf{l_w}}\frac{P(\mathbf{l_w}|\mathbf{x})}{P(\mathbf{l_w})}\right\} \qquad (2)$$

$$\cong \operatorname*{argmax}_{\mathbf{w}}\left\{P(\mathbf{w})\max_{\pi:\pi\in L',\mathcal{B}(\pi_{1:T})=\mathbf{l_w}}\frac{1}{P(\mathbf{l_w})}\prod_{t=1}^{T}y_{\pi_t}^t\right\} \qquad (3)$$
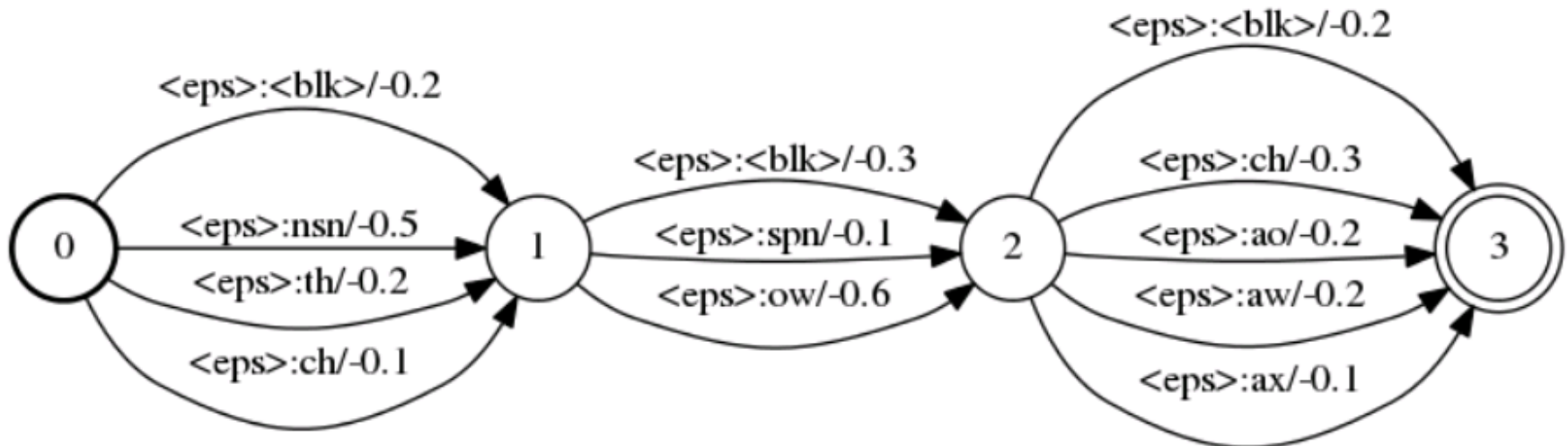
- **frame sync. to phone synchronous decoding**

$$\mathbf{w}^* \cong \operatorname*{argmax}_{\mathbf{w}}\left\{P(\mathbf{w})\max_{\pi:\pi\in L',\mathcal{B}(\pi_{1:T})=\mathbf{l_w}}\frac{1}{P(\mathbf{l_w})}\left\{ \right.\right. \qquad U = \{u : y_{\text{blank}}^u \simeq 1\} \quad (5)$$

$$\left.\left.\prod_{t\notin U}y_{\pi_t}^t \cdot \prod_{t\in U}y_{\text{blank}}^t\right\}\right\} \quad (4)$$

$$= \operatorname*{argmax}_{\mathbf{w}}\left\{P(\mathbf{w})\max_{\pi':\pi'\in L,\mathcal{B}(\pi'_{1:J})=\mathbf{l_w}}\frac{1}{P(\mathbf{l_w})}\prod_{j=1}^{J}y_{\pi'_j}^{t_j}\right\} \quad (6) \qquad J = T - |U| \quad (7)$$

9

# CTC Lattice

- **CTC Lattice - Extremely Compact Acoustic Information Preserver**

| Time | phone label : acoustic score | | | | |
|------|---|---|---|---|---|
| 0.4s | $< blk >: 0.2$ | $nsn : 0.5$ | $th : 0.2$ | $ch : 0.1$ | |
| 0.9s | $< blk >: 0.3$ | $ow : 0.6$ | $spn : 0.1$ | | |
| 1.5s | $< blk >: 0.2$ | $ch : 0.3$ | $ao : 0.2$ | $aw : 0.2$ | $ax : 0.1$ |

# Hypothesis Posterior CM
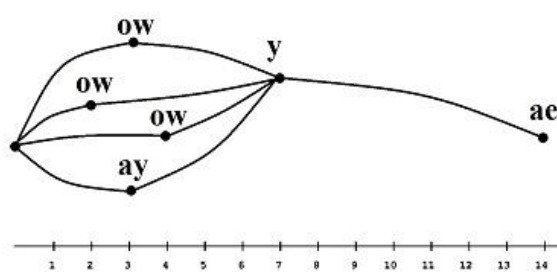
- Procedule:

  - Phone level CTC lattice

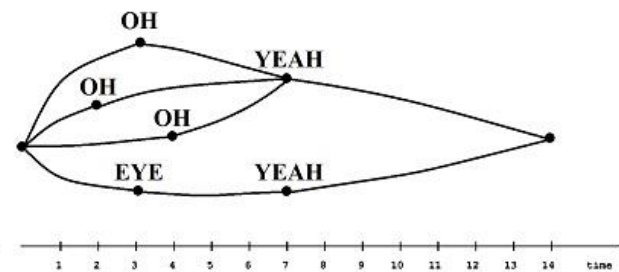  - Word Lattice

  - Confusion Network

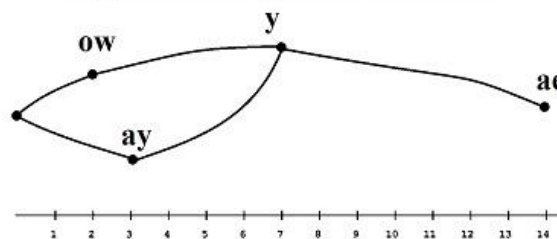  - Confidence Measure



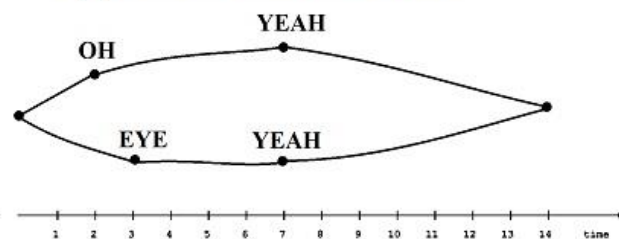(a) Inference results from HMM & CTC

(b) HMM phone lattice

(c) HMM word lattice

(d) CTC phone lattice

(e) CTC word lattice

# Experiment

- Setup
  - Swb 300h, 2-2.5M parameters, NIST hub5e-swb subset
  - details on our paper
- Baseline WER performance

| Model Unit | AM | Decoding | WER |
|------------|-----|----------|------|
| CD-state | DNN-HMM | FSD | 16.7 |
| CI-phone | LSTM-CTC | FSD | 18.7 |
|          |          | PSD | 18.8 |

- CM Evaluation: Normalised Cross Entropy (NCE)

$$NCE = \frac{H(\mathbf{C}) - H(\mathbf{C}|\mathbf{x})}{H(\mathbf{C})}$$

$H(\mathbf{C})$ corresponds to the entropy of the tag sequence,
$H(\mathbf{C}|\mathbf{x})$ is the entropy of the confidence score sequence

  - The higher the better

上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

# Experiment

- Hypothesis Posterior CM [1]

| AM | Decoding | CM | NCE |
|---|---|---|---|
| DNN-HMM | FSD | CN | 0.172 |
| LSTM-CTC | FSD | CN | 0.019 |
| | PSD | CN | 0.224 |
| | | AC+CN | 0.230 |

- CN hypothesis posterior CM can't be directly applied to CI-phone-CTC model
  - Blank allocation problem:
    - e.g., ow \<blk\> ch \<blk\> \<blk\> \<blk\> ao \<blk\>

---

[1] We also derive a PSD version of predictor based CM, detail comparison can be referred to our paper.

13

上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

# Experiment

- Hypothesis Posterior CM [1]

| AM | Decoding | CM | NCE |
|---|---|---|---|
| DNN-HMM | FSD | CN | 0.172 |
| LSTM-CTC | FSD | CN | 0.019 |
| | PSD | CN | 0.224 |
| | | AC+CN | 0.230 |

- CN hypothesis posterior CM can't be directly applied to CI-phone-CTC model
  - Blank allocation problem:
    - e.g., ow <blk> ch <blk> <blk> <blk> ao <blk>
- In PSD, CN hypothesis posterior CM can be successfully applied
- Even with significantly better NCE: 0.224 → 0.172

[1] We also derive a PSD version of predictor based CM, detail comparison can be referred to our paper.

上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

- Reason of Better CM
  - Better lattice
    - Phone



Oracle phone error rate

OPER (%)

Legend: FSD CTC, PSD CTC, FSD HMM

Phone Lattice Density (arc/frame)
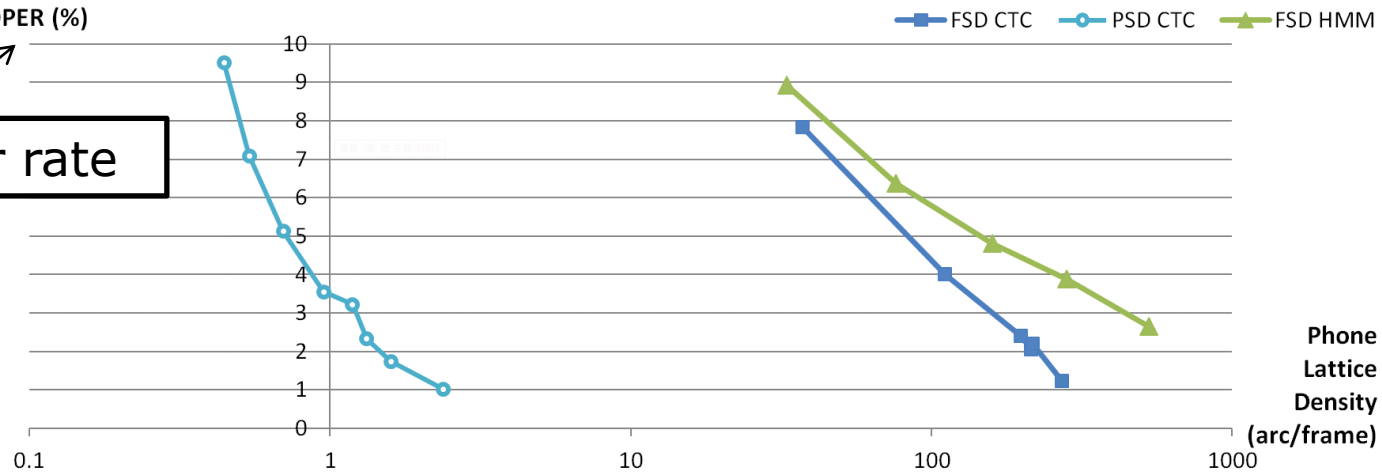
上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

# Experiment

- **Reason of Better CM**
  - **Better lattice**

Phone

OPER (%)

Oracle phone error rate

Word

$$1 - \frac{OWER}{WER}$$



Phone Lattice Density (arc/frame)

Relative Oracle Word Error Rate Reduction (%)

Word Lattice Density (arc/frame)

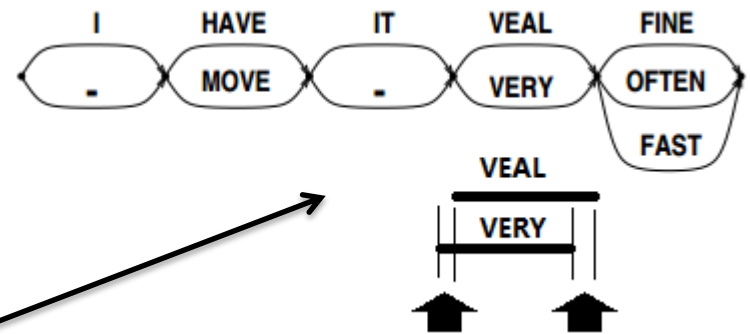# Experiment

- Reason of Better CM
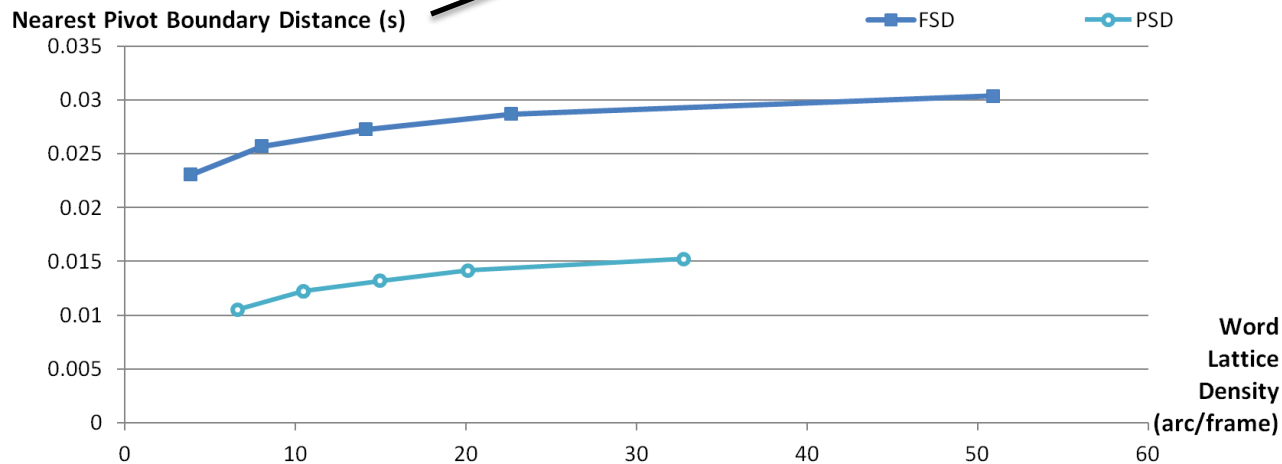  - Larger CN depth → more competing information

# Experiment

- Reason of Better CM
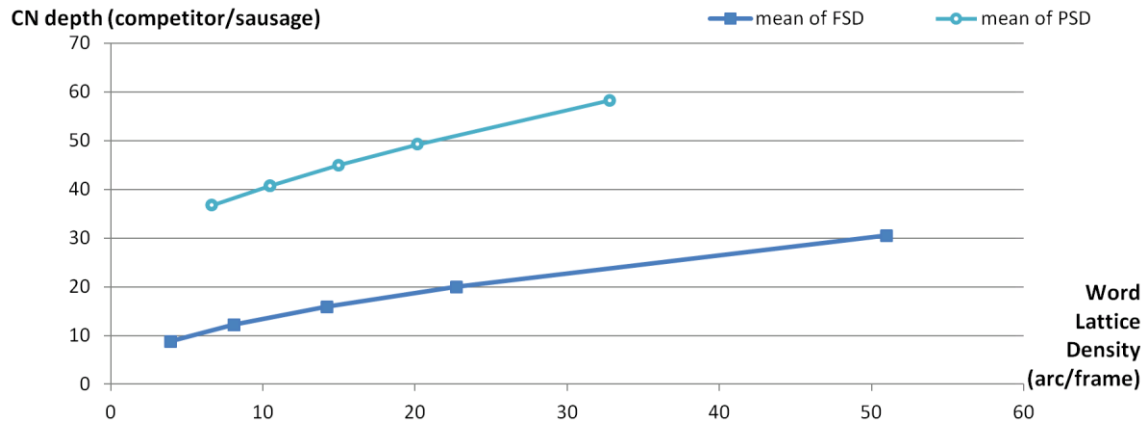  - Larger CN depth → more competing information
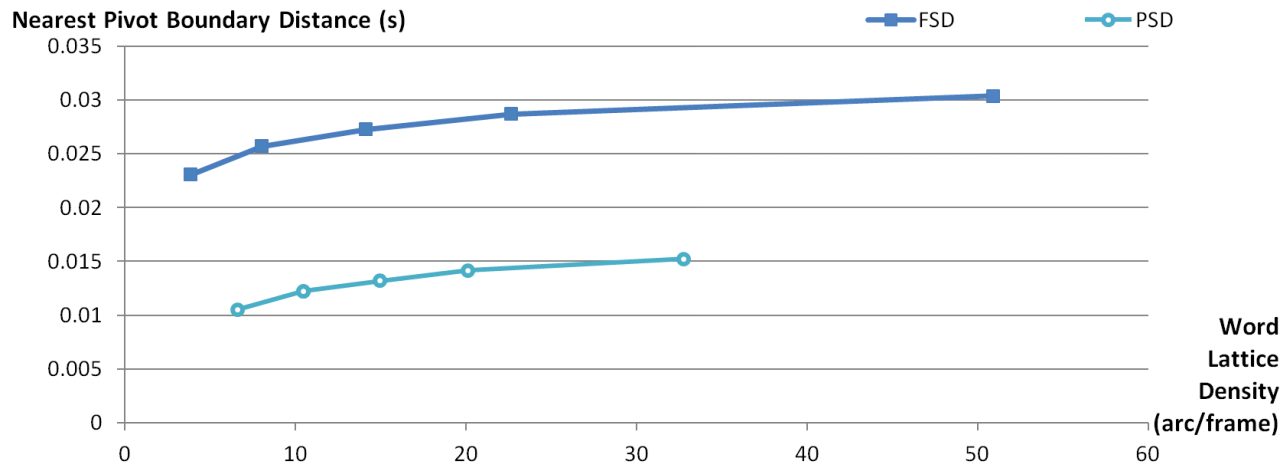


  - Because of more stable boundary



18

# Experiment

- Reason of Better CM
  - Larger CN depth → more competing information



  - Because of more stable boundary

# Summary

- The potential of **compact** and **precise** PSD CTC lattice in preserving acoustic information was utilized to form better CMs
- PSD version of predictor based CM was proposed with elaborate phonemic normalization and blank info (in paper)
- The characteristics of **lattice** and **confusion network** generated from **PSD** framework were carefully investigated, and CN hypothesis posterior CM was proposed
- The two types of CMs can be combined together as a pair of complements

- Future work: applying proposed CMs as predictors in model training framework

上海交通大学
SHANGHAI JIAO TONG UNIVERSITY