# Sphinx Benchmark Report

Long Qin

Language Technologies Institute School of Computer Science Carnegie Mellon University

#### Overview

- O Evaluate general training and testing schemes
  - O LDA-MLLT, VTLN, MMI, SAT, MLLR, CMLLR
- O Use default setup and existing tools
  - O SphinxTrain-0.8, Sphinx3
- O Focus on WER, running time was not measured
  - Experiments were performed on different server machines, it's not easy to directly compare the xRT
- O Test on different data
  - O Easy task (WSJ) vs. broadcast news
  - English vs. Mandarin

## Outline

- O The baseline training scheme
- O LDA-MLLT
- O VTLN
- o MMI
- O SAT
- O CMLLR
- O MLLR
- O Experiments
- O Discussion

## Baseline Training Scheme

13-MFCC with Delta and Delta-Delta

Triphone model
3-state HMM
GMM observation distribution

Feature Extraction



CI Model



CD Model

Monophone model
3-state HMM
1-Gaussian or GMM
observation distribution

Decision tree clustering with auto-generated questions
A few thousand states

## Force Alignment



- Force Alignment
  - Find the best alignment between speech and corresponding HMMs
- O Goal
  - Possibly remove utterances with transcription errors or low quality recordings
  - Find appropriate pronunciations for words with multiple pronunciations
- Settings
  - \$ \$CFG\_FORCEDALIGN = "yes";
  - \$CFG\_FORCE\_ALIGN\_BEAM = 1e-60;
  - \$CFG\_FALIGN\_CI\_MGAU = "yes"/"no";

### LDA-MLLT



- O LDA (linear discriminant analysis)
  - O Find a linear transform of feature vector, so that class separation is maximized
  - Reduce feature dimension
- MLLT (maximum likelihood linear transform)
  - Minimize the loss of likelihood between full and diagonal variance model
  - Applied together with LDA
- O In Sphinx
  - O Each Gaussian is considered as one class
    - Easier to implement
    - O Could also define state or phone as class
- Settings:
  - \$CFG\_LDA\_MLLT = "yes";
  - \$CFG\_LDA\_DIMENSION = 29;

#### VTLN



- O VTLN (vocal tract length normalization)
  - Formant frequency is considered to have a linear relationship with the vocal tract length
  - Adjust vocal tract length for each speaker to an average length by warping their spectra
  - The warping factor:

$$\lambda = \operatorname{arg\,max} P(O \mid X, \lambda_k)$$

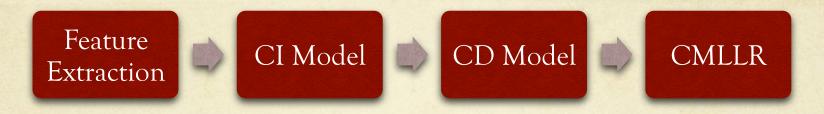
- O In Sphinx
  - O Warping factor is estimated for each utterance using exhaust search
    - O Could also estimate identical warping factor for each speaker
  - O Warping factor should be estimated in both training and decoding
- O Settings:
  - \$CFG\_VTLN = 'yes';
  - \$CFG\_VTLN\_START = 0.70;
  - \$CFG VTLN END = 1.40;
  - \$CFG VTLN STEP = 0.05;

#### MMI



- MMI (maximum mutual information)
  - A discriminative training algorithm
  - Maximize the posterior probability of the true hypothesis
  - Training is time consuming
- O Settings:
  - \$CFG\_MMIE\_MAX\_ITERATIONS = 4;
  - \$CFG\_MMIE\_CONSTE = "3.0";
  - \$CFG\_LANGUAGEWEIGHT = "11.5";
    - O The same as the language weight used in decoding
  - \$CFG\_LANGUAGEMODEL = "LMFILE";
    - O A unigram or bigram LM

#### CMLLR



- O CMLLR (constraint maximum likelihood linear regression)
  - A speaker adaptation algorithm to modify speaker independent system towards new speaker using limited data
  - O Use the same transform for both mean and variance, therefore usually require less data then MLLR
  - Could be formulated as a linear transform of input features

#### O In Sphinx

- O Use a single global transform to adapt the input features for each speaker
- When accumulate counts, run BW with "-fullvar yes", "-2passvar no" and "-cmllrdump yes"

#### O Settings:

- \$CFG\_DEC\_DICTIONARY = "DECODING\_DICTIONRY";
- \$CFG\_DEC\_LM = "DECODING\_LANGUAGE\_MODEL";

#### SAT



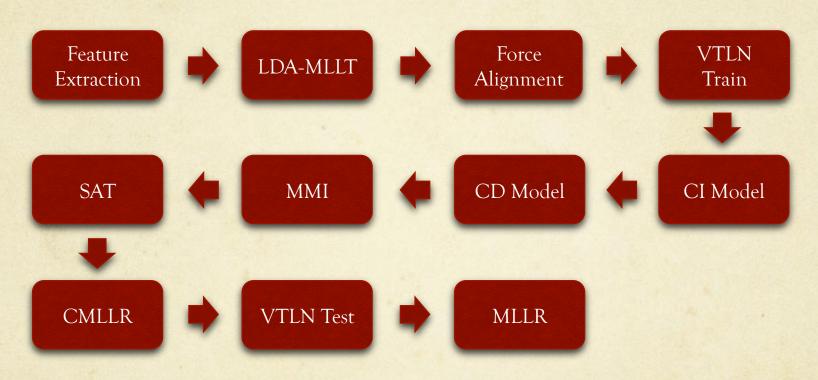
- O SAT (speaker adaptive training)
  - O Train a better speaker independent system
  - Apply CMLLR transforms to training features
  - Re-estimate the CMLLR transforms every iteration
- O In Sphinx
  - O SAT is applied after training a fairly good ML/MMI model
  - Need to split the training control and reference files into smaller files for each speaker (make\_speaker\_lists.py)
- O Settings:
  - \$CFG\_SAT\_DIR = "\$CFG\_BASE\_DIR/sat";

#### MLLR



- O MLLR (maximum likelihood linear regression)
  - Another speaker adaptation algorithm
  - Adjust mean and/or covariance to maximize the likelihood of the adaptation data
- O In Sphinx
  - Adapt mean in default
    - O Could also adapt covariance
  - O Use a single global transform for all models
    - O Could have multiple transforms for different classes of models
- Settings
  - Applied during decoding
    - O Get hypotheses of the testing data from the first pass decoding
    - O Using those hypotheses and testing data to estimate transforms and update model parameters
      - O During bw run, must set "-2passvar no"
    - O Decode again using the adapted model
  - O It's the same procedure when we apply CMLLR/VTLN in decoding

# Overall System Framework



# Data

Training		Testing	LM
WSJ0	15-hour	Nov. 92 5k and 20k Dev/Eval	standard Trigram
WSJ0+1	82-hour		
BN	138-hour	HUB4-96 Dev/Eval (with data in all different environments)	Trigram from BN 92-97 LM data
Mandarin BN	128-hour	RT04-Eval	Trigram from Chinese Gigaword

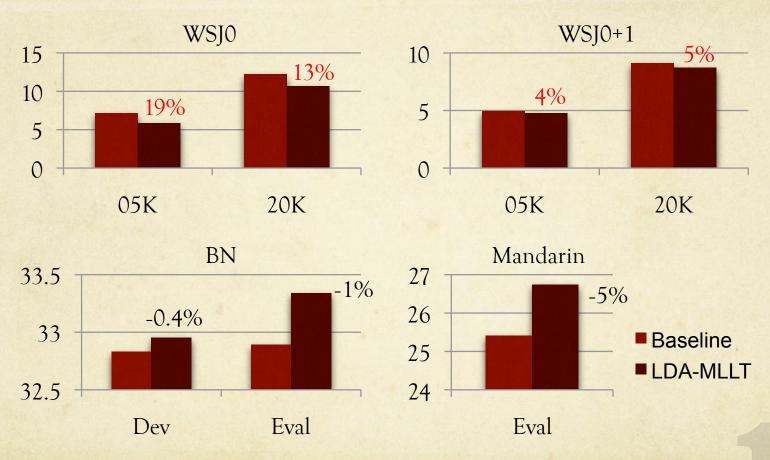
# Baseline Settings

- O Force Alignment
  - O Could use multiple-Gaussian CI model
    - A little bit better, more computation
- C Linguistic Questions
  - If available
  - Or use auto-generated questions
- O Decoding
  - 0 lw=11.5, beam=1e-100, wbeam=1e-80, wip=0.2
- Mixtures and States
  - O WSJ0: 16 mixtures, 2000 tied-states
  - O WSJ0+1: 32 mixtures, 4000 tied-states
  - O BN: 32 mixtures, 5000 tied-states
  - Mandarin: 32 mixtures, 4000 tied-states

## Baseline Results

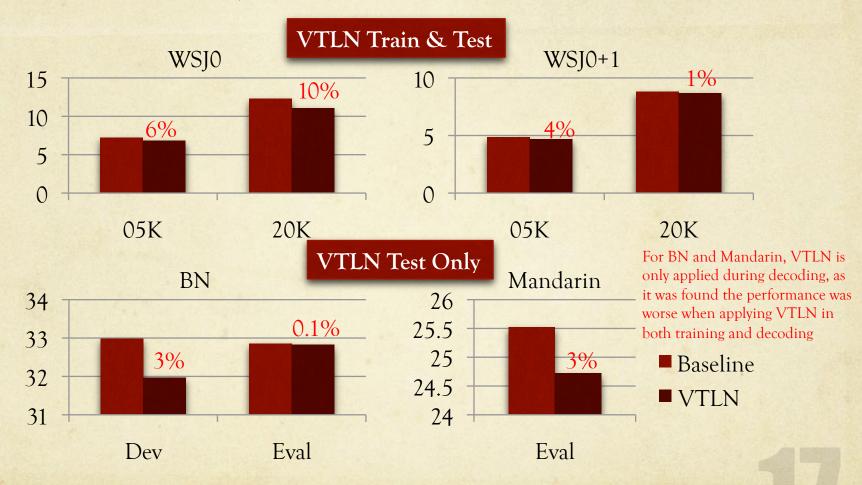
	WER (%)		
Data	Dev	Eval	
WSJO	7.62 (5k) 12.84 (20k)	6.85 (5k) 11.69 (20k)	
WSJ0+1	5.50 (5k) 9.80 (20k)	4.18 (5k) 7.78 (20k)	
BN	32.98	32.85	
Mandarin		25.35	

## LDA-MLLT Results



Comment: may work better on simple tasks with high quality data, but others (Joao Miranda) had tried it on noisy data, which also helped a lot. It works on telephone conversation tasks too.

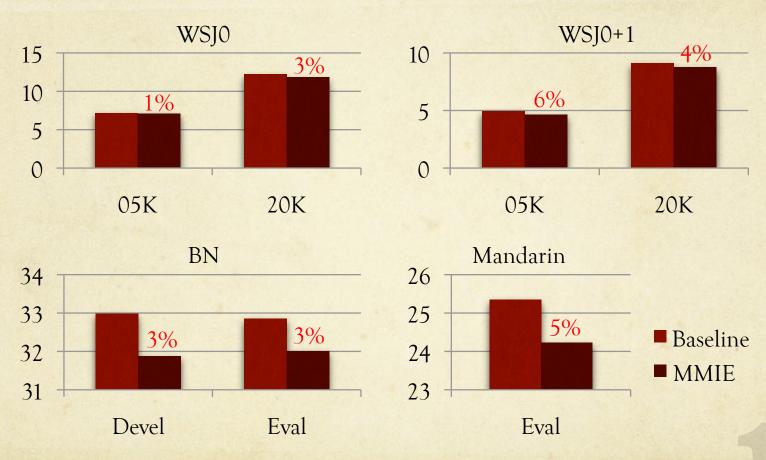
## VTLN Results



#### To be noticed:

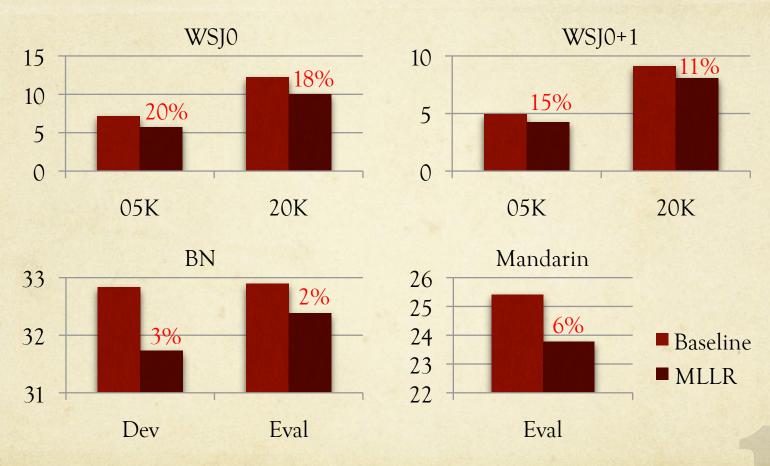
- the red numbers in the graph is the relative improvement over the baseline
- to have a graph without too many bars, the WSJ 5K/20K results are the average of the the Dev and Eval results

#### MMI Results



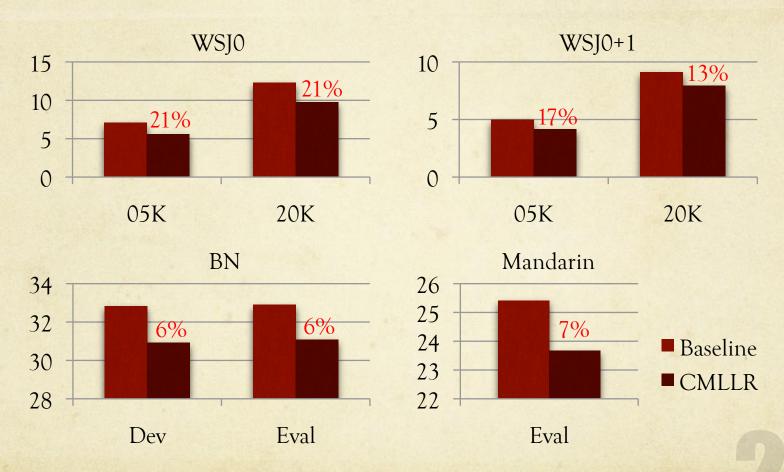
Comment: the results are not as good as I got from the lattice pruning experiments, where I used smaller lattices; try smaller beam widths when generate lattices, such as \$beam = \$wbeam = 1e-70, should be better and faster. Also try to use a bigram instead of unigram when generating lattices.

## MLLR Results



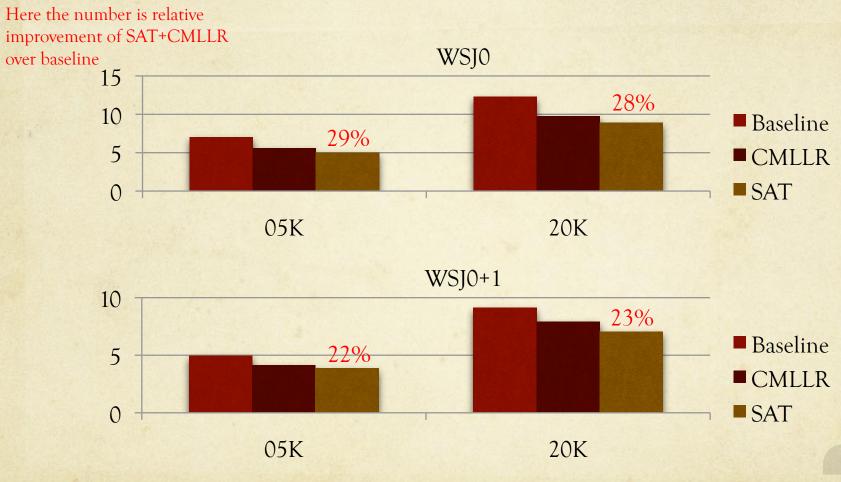
Comment: works pretty good especially when the first path hypotheses are accurate; could use the second path hypotheses train a better transform and iteratively do it to get the best number

## CMLLR Results



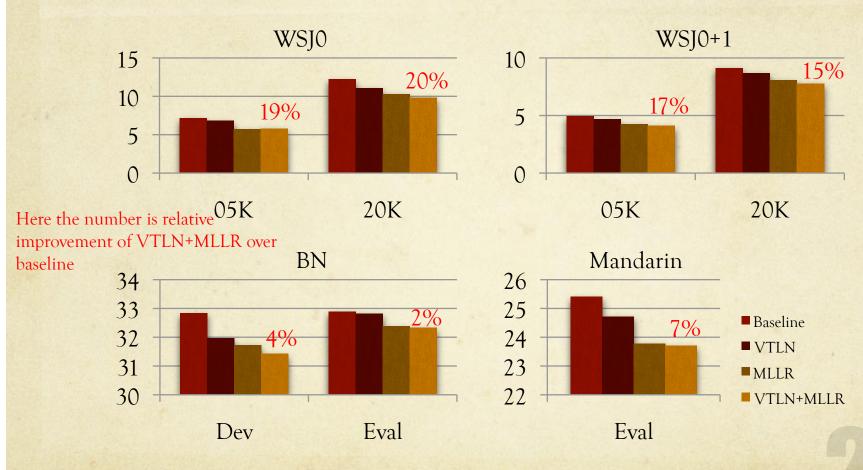
Comment: has similar performance as MLLR, slightly better in BN

### SAT Results



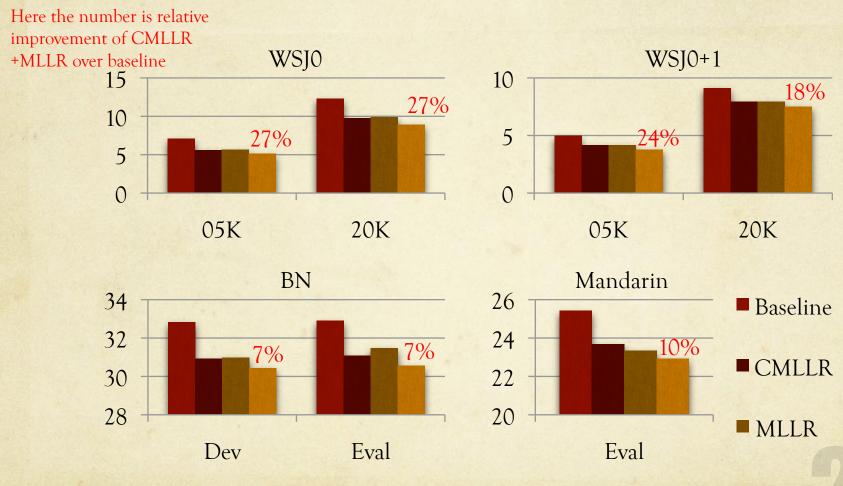
Comment: SAT + CMLLR decoding is very effective, which usually gives 10% improvement over CMLLR decoding only. When estimating CMLLR transform, it's better to start from a very good hypothesis such as the CMLLR+MLLR decoding result.

### VTLN + MLLR Results



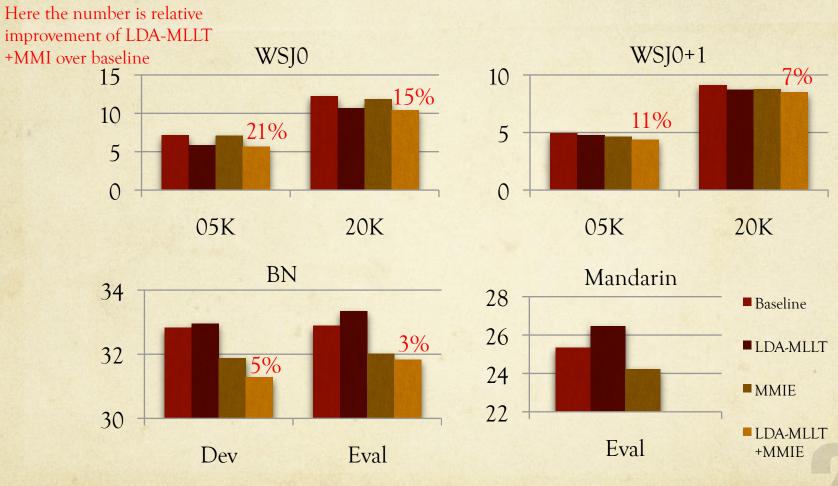
Comment: the improvement is additive, but quite small compared to perform MLLR only

## CMLLR+MLLR Results



Comment: CMLLR+MLLR further improves the WER!

## LDA-MLLT + MMI Result



Comment: MMIE gives solid improvement over LDA-MLLT (compare the 2<sup>nd</sup> bar and the 4<sup>th</sup> bar)

# Summary

- O LDA-MLLT
  - o works pretty good on simple tasks with clear speech, not clear on hard tasks with noisy speech, needs more investigation
- O VTLN
  - o produces some improvement
- O MMIE
  - o produces ok/good improvement
  - o requires large amount computation
- O CMLLR
  - o works pretty good, especially when first path hypotheses are very accurate
- O MLLR
  - o works similar to CMLLR
- O SAT
  - o produces solid improvement

# Still Missing

- O Better discriminative training technique
  - O boosted-MMI
- O Deep Neutral Network
  - O Bottle Neck Feature (easier to adapt)
  - O Hybrid Model (more improvement)