Speech Recognition Techniques for a Sign Language Recognition System

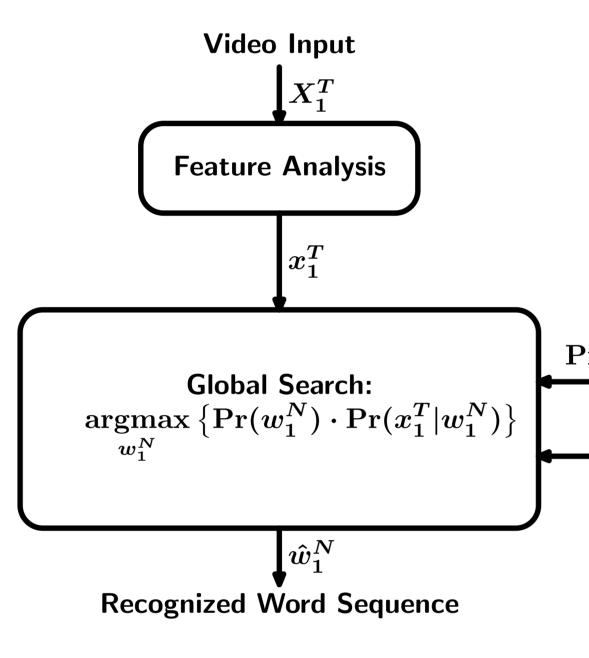
Introduction

- automatic sign language recognition system
- necessary for communication between deaf and hearing people
- continuous sign language recognition, several speakers, vision-based approach, no special hardware
- Iarge vocabulary speech recognition (LVSR) system to obtain a textual representation of the signed sentences
- evaluation of speech recognition techniques on publicly available sign language corpus

Automatic Sign Language Recognition (ASLR)

- similar to speech recognition: temporal sequences of images
- important features
- hand-shapes, facial expressions, lip-patterns
- orientation and movement of the hands, arms or body
- HMMs are used to compensate time and amplitude variations of the signers

goal: find the model which best expresses the observation sequence



Experimental Setup

Database

- system evaluation on the RWTH-BOSTON-104 database
- ▶ 201 sentences (161 training and 40 test sequences)
- vocabulary size of 104 words
- Speakers (2 female, 1 male)
- corpus is annotated in glosses

Problems

- 26% of the training data are singletons
- simple sentence structure
- one out-of-vocabulary (OOV) words with whole-word models

Differences in Comparison to ASR

- simultaneousness
- signing space
- environment
- speakers and dialects
- coarticulation and movement epenthesis
- ▶ silence
- whole-word models and sub-word units

Philippe Dreuw, David Rybach, Thomas Deselaers, Morteza Zahedi, and Hermann Ney

Human Language Technology and Pattern Recognition, RWTH Aachen University, Aachen, Germany

| $\mathrm{r}(x_1^T w_1^N)$ | Word Model Inventory |
|---------------------------|----------------------|
| $\Pr(w_1^N)$ | Language Model |









System Overview

Visual Modeling (VM)

- related to the acoustic model in ASR
- HMM based, with separate GMMs, globally pooled diag. covariance matrix
- monophone whole-word models
- pronunciation handling

Language Modeling (LM)

- according to ASR: LM should have a greater weight than the VM
- trigram LM using the SRILM toolkit, with modified Kneser-Ney discounting with interpolation

Features

- appearance-based image features: for baseline system
- thumbnails of video sequence frames (intensity images scaled to 32x32 pixels)
- give a global description of all (manual and non-manual) features proposed in linguistic research
- manual features:
- dominant hand tracking: hand position, hand velocity, and hand trajectory features

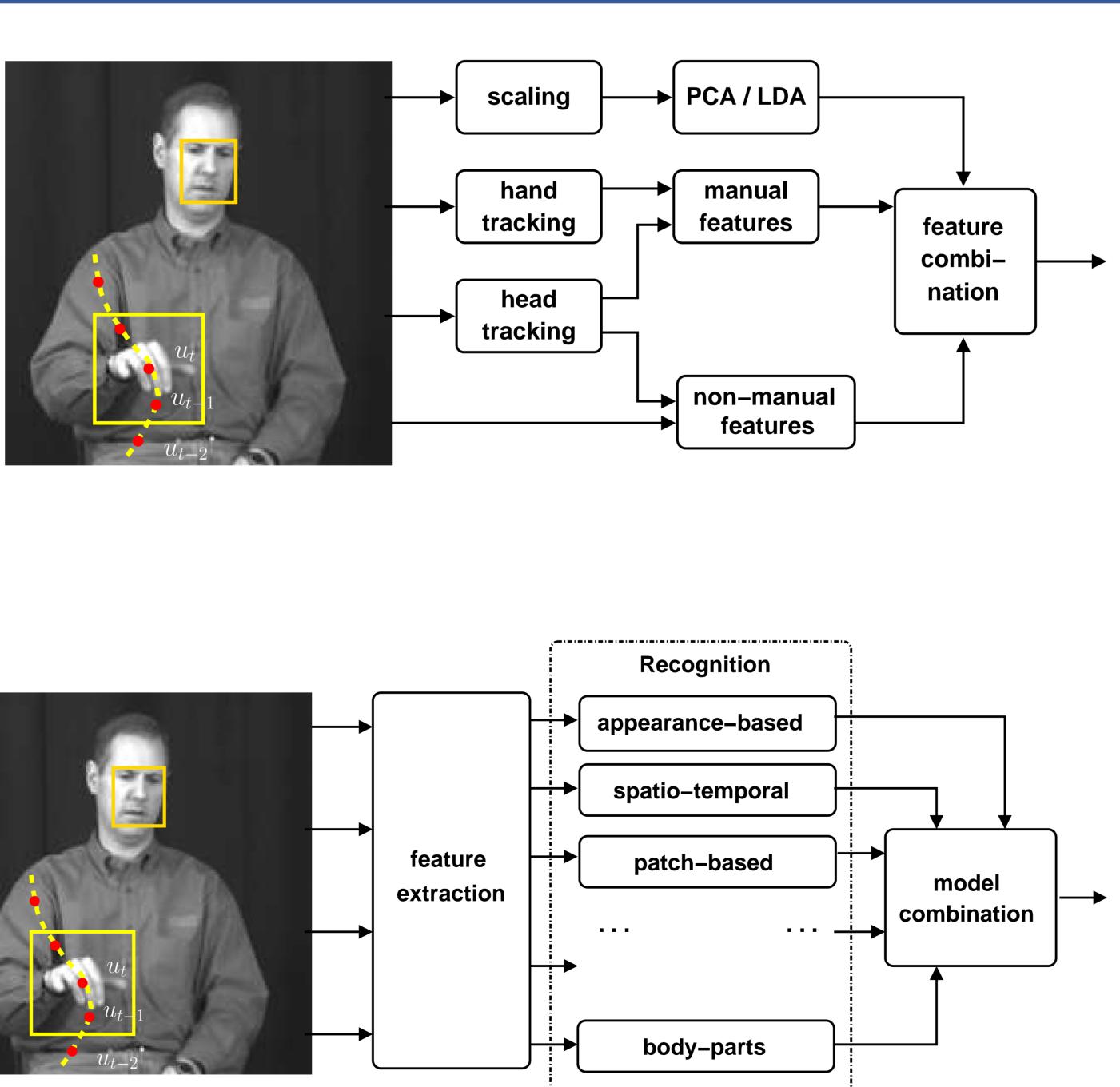
Feature Selection and Model Combination

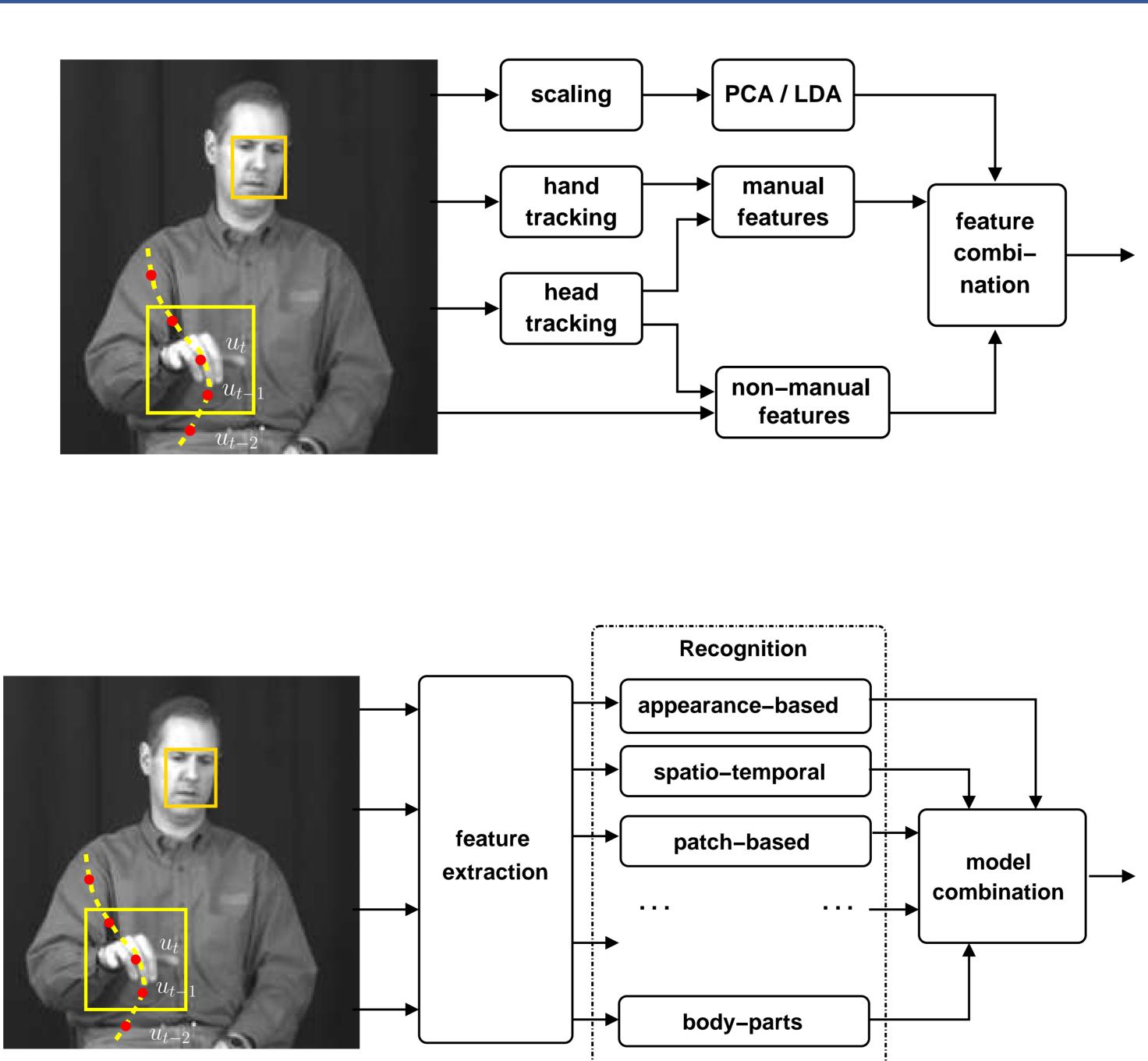
Feature Selection

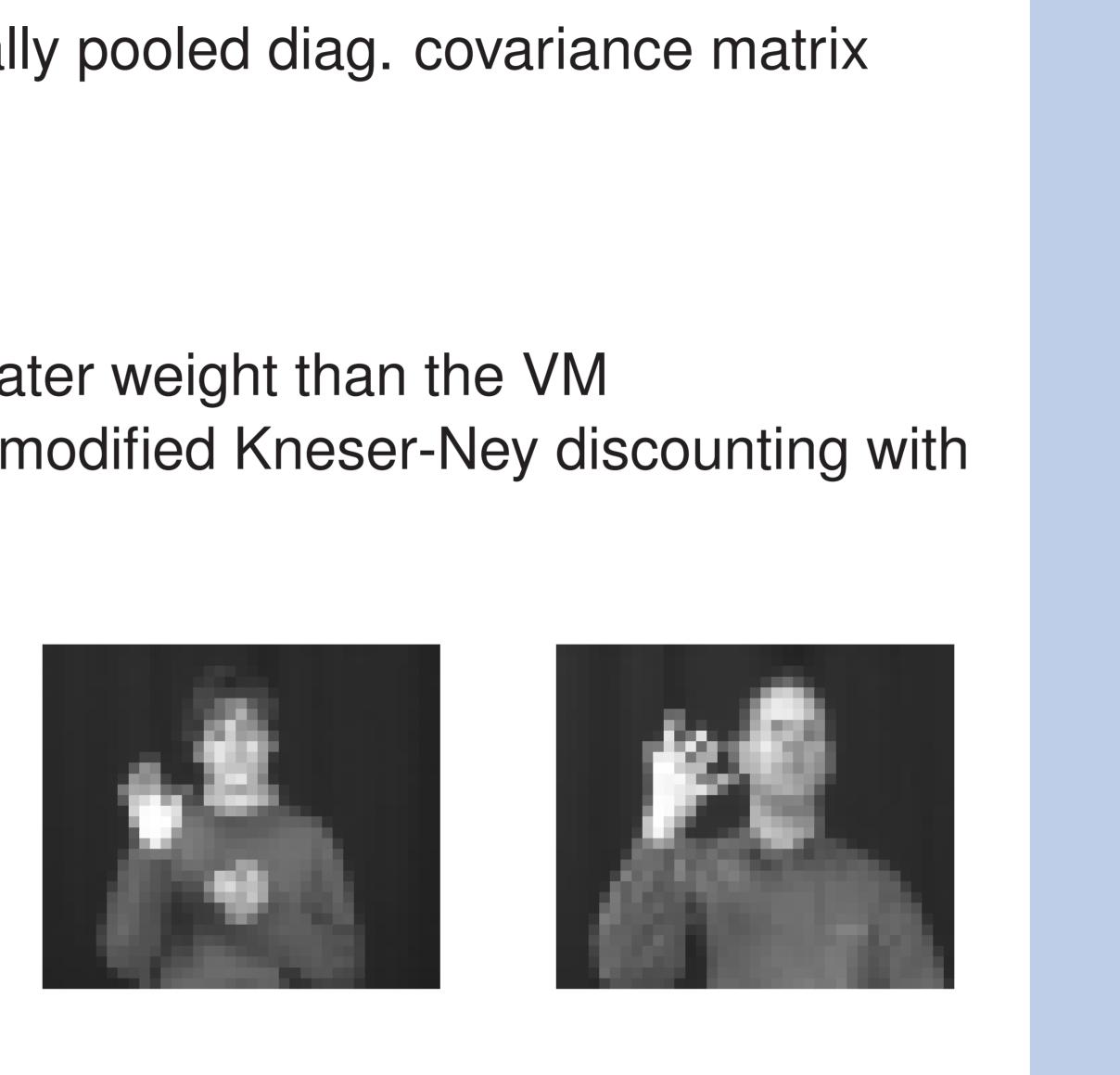
- concatenation of appearance-based and manual features
- sliding window for context modeling
- dimensionality reduction by PCA and/or LDA

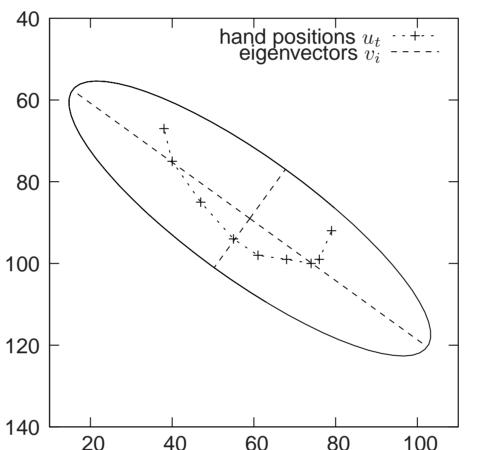
Model Combination

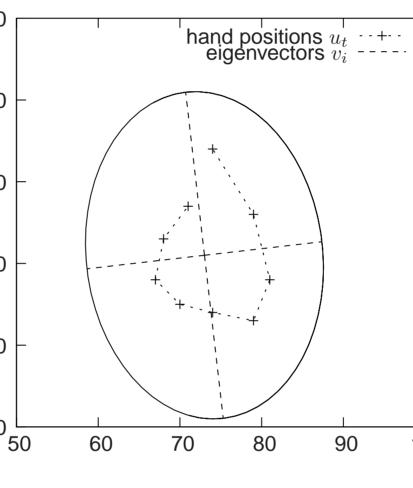
- Iog-linear combination of independently trained models
- profit from independent alignments (e.g. performing well for long and short words)
- profit from different feature extraction approaches

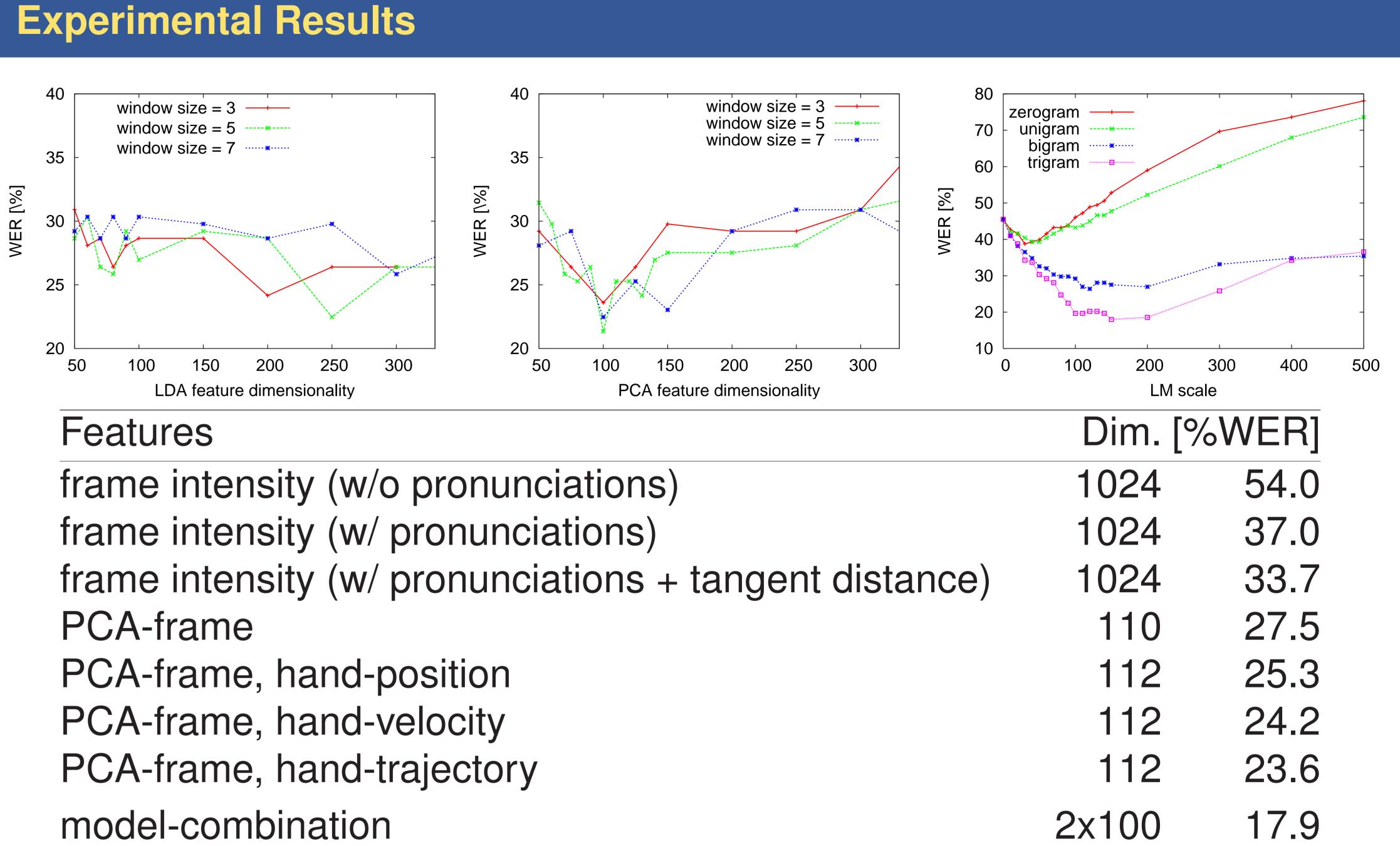


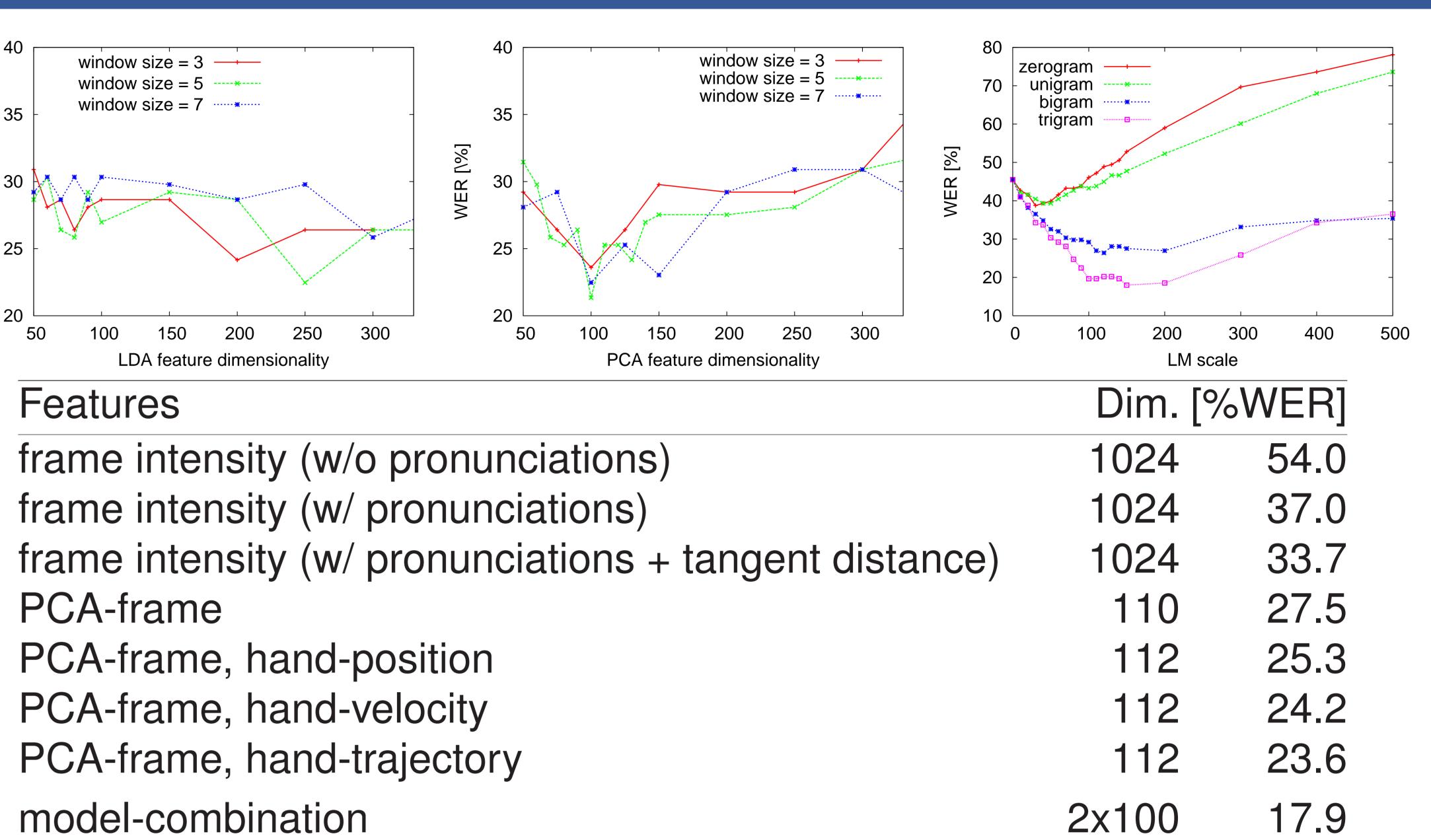












Example Results

Correct Examples IX-1P FIND SOMETHING-ONE BOOK IX-1P FIND SOMETHING-ONE BOOK LOVE JOHN WHO LOVE JOHN WHO JOHN BUY YESTERDAY WHAT BOOK JOHN BUY YESTERDAY WHAT BOOK

Incorrect Examples

| MARY | VEGETABLE KI |
|------|--------------|
| MARY | VEGETABLE KI |
| JOHN | IX GIVE MA |
| JOHN | IX WOMAN |
| | LIKE CHOCOLA |
| JOHN | LIKE CHOCOLA |
| JOHN | [UNKNOWN] |
| JOHN | FUTURE NO |
| | |

Conclusion

- modelling, and model combination
- outlook: connection of recognizer output to a statistical machine translation system achieved promising translation results



RWTH-BOSTON-104 Database

Corpus Statistics Training Test JOHN FISH WONT EAT BUT CAN EAT CHICKEN 161 sentences 40 JOHN FISH WONT EAT BUT CAN EAT CHICKEN 710 178 running words 12422 3324 frames vocabulary 103 65 27 singletons OOV **LM Perplexities** KNOW IX LIKE CORN PPLM type KNOW IX LIKE MARY 106.0 zerogram AN IX NEW COAT 36.8 unigram ____ NEW COAT 6.7 bigram ATE WHO 4.7 trigram ATE WHO

BUY HOUSE Database is publicly available OT BUY HOUSE

LVSR system is suitable for vision-based continuous sign language recognition many of the principles known from ASR can directly be transfered Important for ASLR: temporal contexts, pronunciation handling, language