# Building language detectors with small amounts of data

David van Leeuwen (TNO)
Niko Brümmer (SDV)

TNO | Knowledge for business

Leaders in voice Transaction Management

Spescom DataVoice
Leaders in Voice Transaction Management

SPESCOM

# Synopsis

- Standard language model training for language recognition needs lots of data
  - typically 60 hours speech, 100 speakers, per language

- We would like to reduce this demand

- Investigate classifier that works in *score space* rather than acoustic space

- Evaluate with
  - LRE-2005 (7 languages)
  - CSLU-22 (21 languages)

- Can train score-space based system with ~ 1 hour data
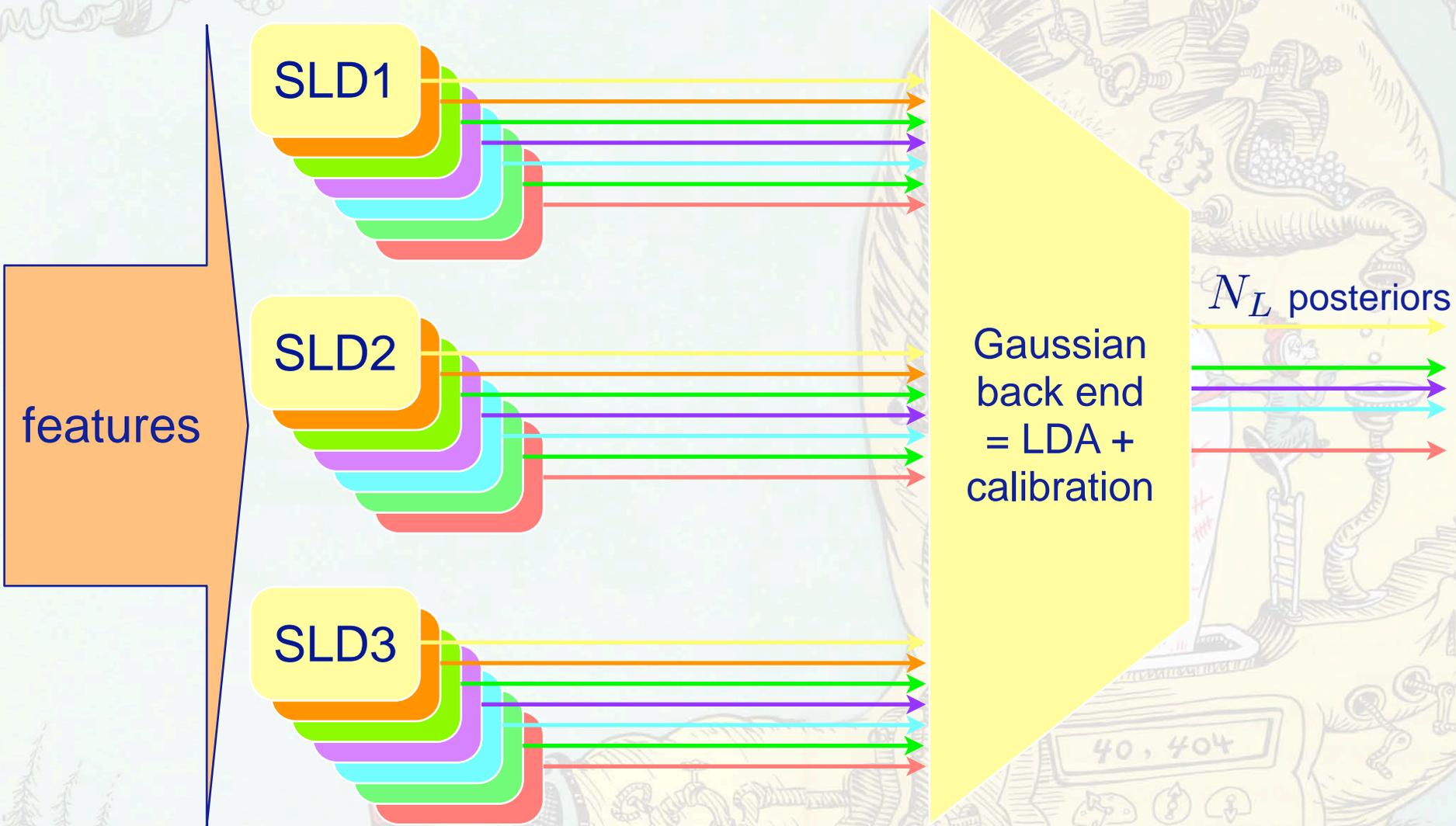  - at twice the $C_{DET}$

# Motivation

- Sometimes there is not much training data for new language available
    - e.g., Indian accented English in LRE-2005: 20 minutes

- Sometimes we may not *want* to train acoustic model for new language
    - Hard for inexperienced user

- Notion that language recognition *back-end* can 'repair' sub-optimal modeling performance
    - Try to let back-end to the whole job, without specific acoustic language model

# Caveat

- Collection of large amounts of speech *should* be relatively easy
    - no orthographic annotation required

- But:
    - correct labeling of language *is* required
    - different collection characteristic to background data will lead to confounding of *language* and *data collection* modeling

- This is true for *any* kind of modeling
    - front-end (GMM, SVM, acoustics, phonotactics)
    - back-end (LDA, logistic regression)

LDA: Linear Discriminant Analysis
SVM: Support Vector Machine
GMM: Gaussian Mixture Model

# Overview of (typical) LID system

features

SLD1

SLD2

$N_L$ posteriors

Gaussian
back end
= LDA +
calibration

SLD3

SLD: Single Language Detector

# Modeling power of LDA back-end

- with proper priors and threshold for posterior
  - optimal NIST LRE decisions can be made
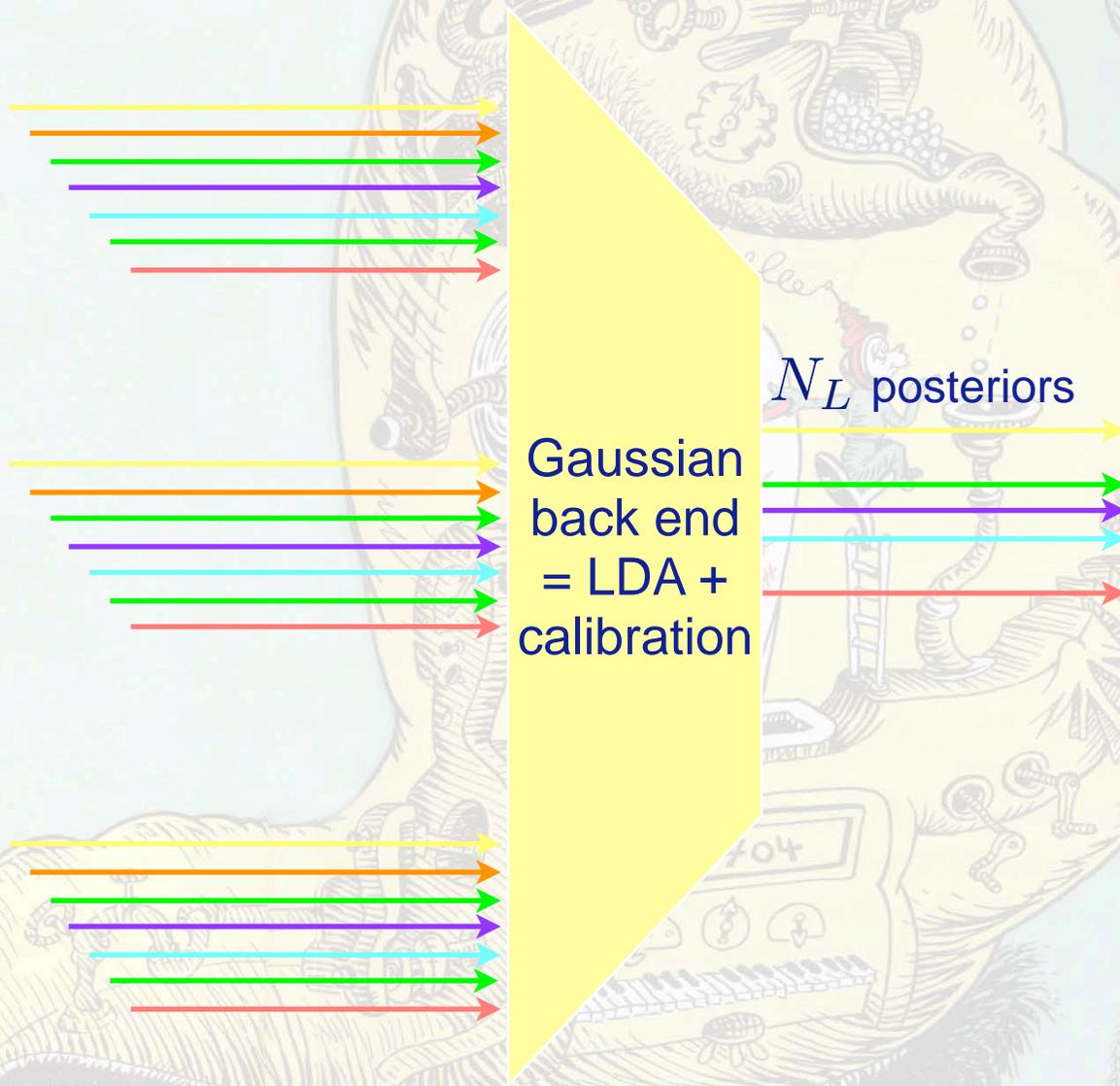
- LRE-2005 languages

$$p(L_{2005}) = 1/N_L$$

- Other CallFriend

$$p(L_{\text{CF}\backslash 2005}) = 0$$
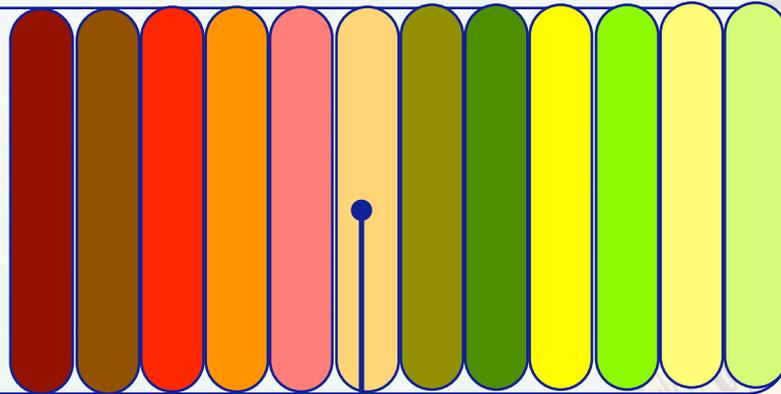
- Posterior threshold

$$\theta = 1/N_L$$

Gaussian back end = LDA + calibration

$N_L$ posteriors

# LRE-2005: Jack-knifing approach

- for each target language $L_i$

    - remove Single Language Detector $L_i$ from LDA training

    - build LDA, using all LDA training trials (incl. $L_i$)

    - compute target and non-target scores for these $L_i$ test-segments, and make decisions

- pool decisions, calculate $C_{DET}$ according to NIST LRE plan

LDA: Linear Discriminant Analysis
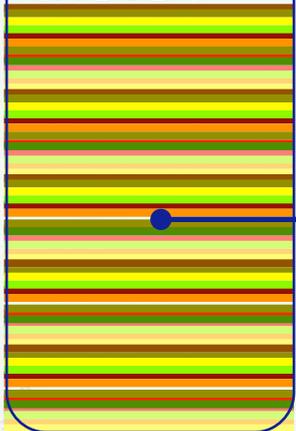$C_{DET}$: Cost of detection
LRE: Language Recognition Evaluation

# Application to NIST LRE-2005

**CallFriend training 12 languages**

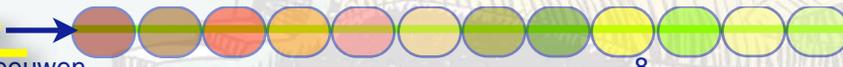**Single Language Detectors**

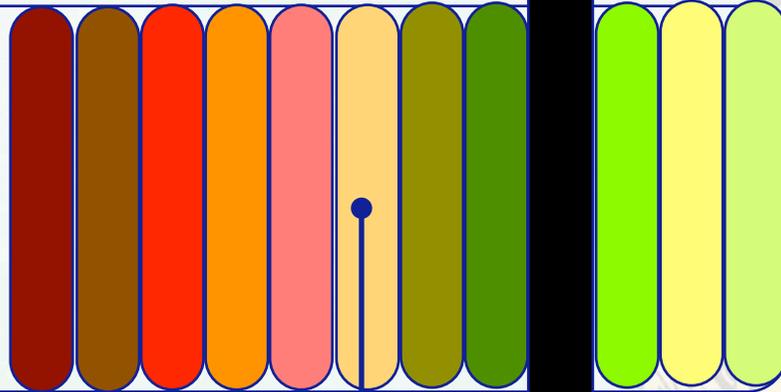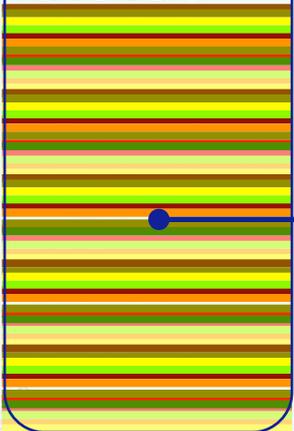**LDA training trials**

**Linear Discriminant Analysis training**

**LDA Back-end**

Test trials

Score

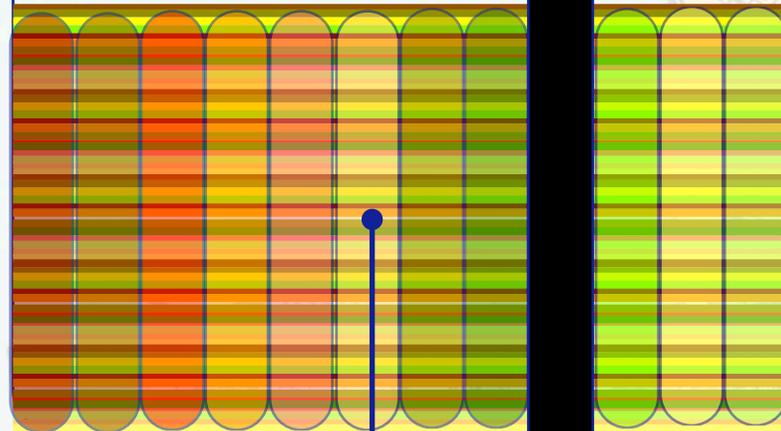# Application to NIST LRE-2005

**CallFriend training 12 languages**

**LDA training trials**

**Linear Discriminant Analysis training**

**Single Language Detectors**

**LDA Back-end**

Test trials

Score

# Application to NIST LRE-2005

**CallFriend training 12 languages**

**Single Language Detectors**

**LDA training trials**
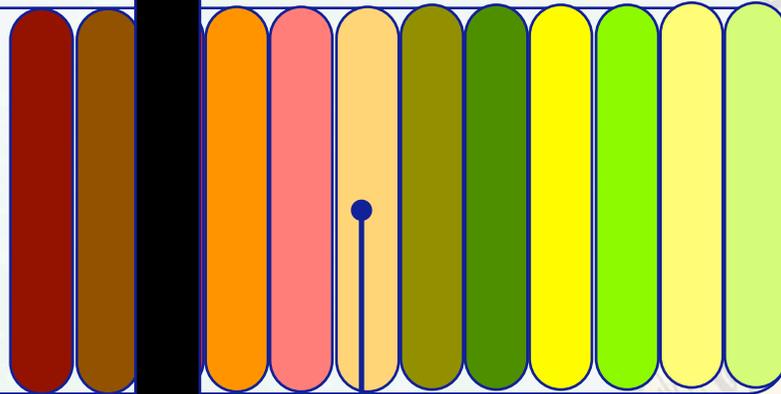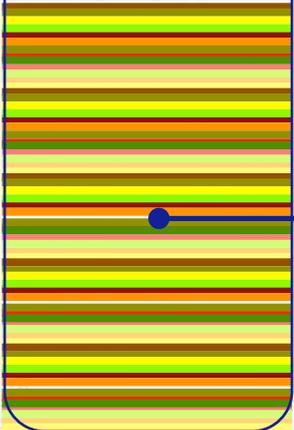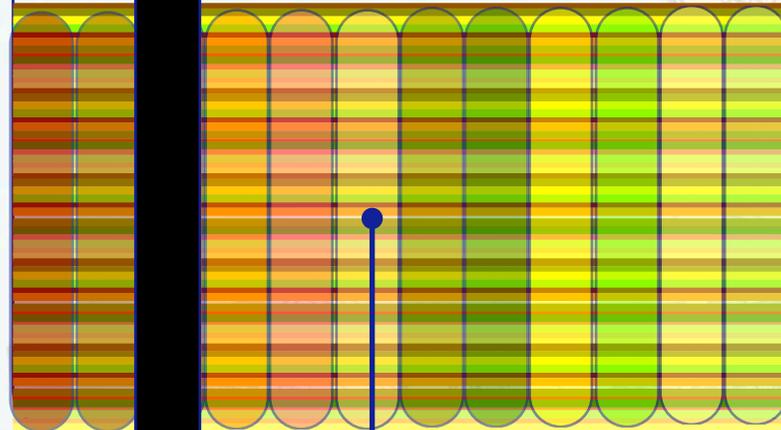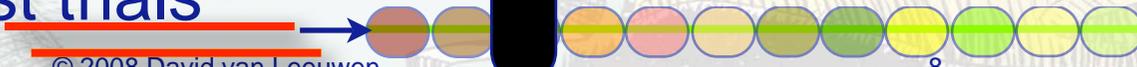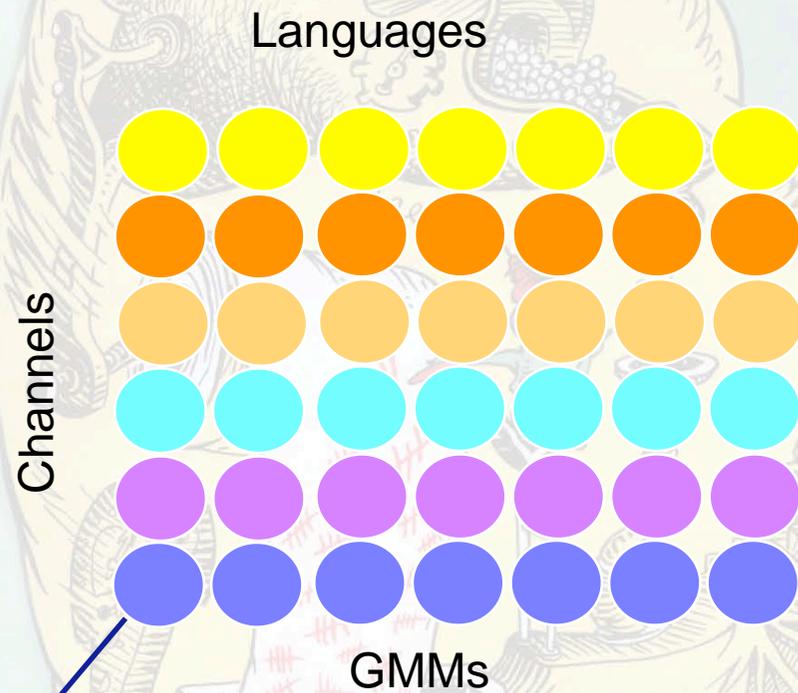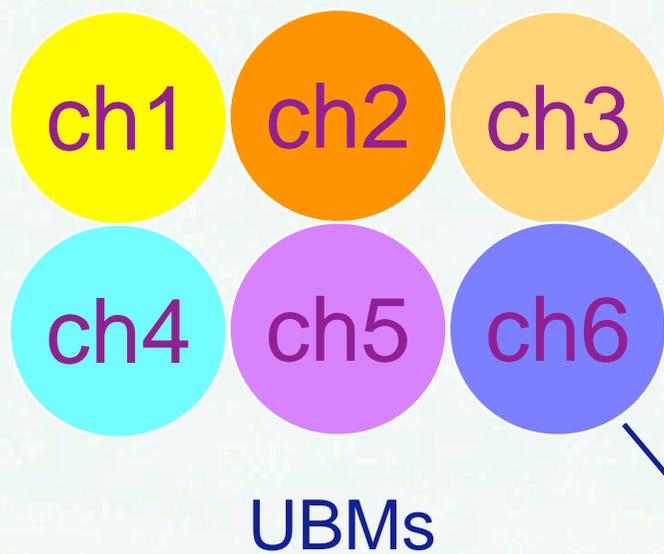
**Linear Discriminant Analysis training**

**LDA Back-end**

Test trials

Score

8

# Three systems: 1) Chan-GMM

ch1 ch2 ch3
ch4 ch5 ch6

UBMs

Languages

Channels

GMMs

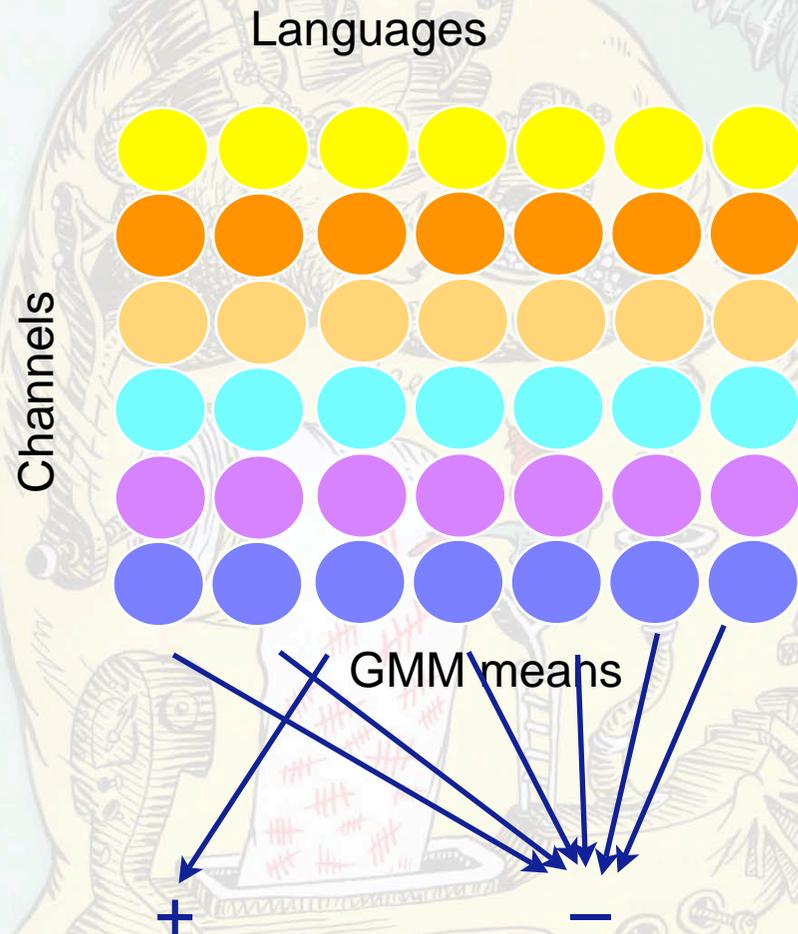$N_{ch} \times N_L$
Likelihood ratios

# Three systems: 2) GMS (GMM means SVM)

s1

s2

Sexes

Languages

GMM means

+

−

SVM

$2\,N_L$
Scores

10

# Three systems: 3) Chan-GMS

ch1 ch2 ch3
ch4 ch5 ch6

UBMs

Languages

Channels
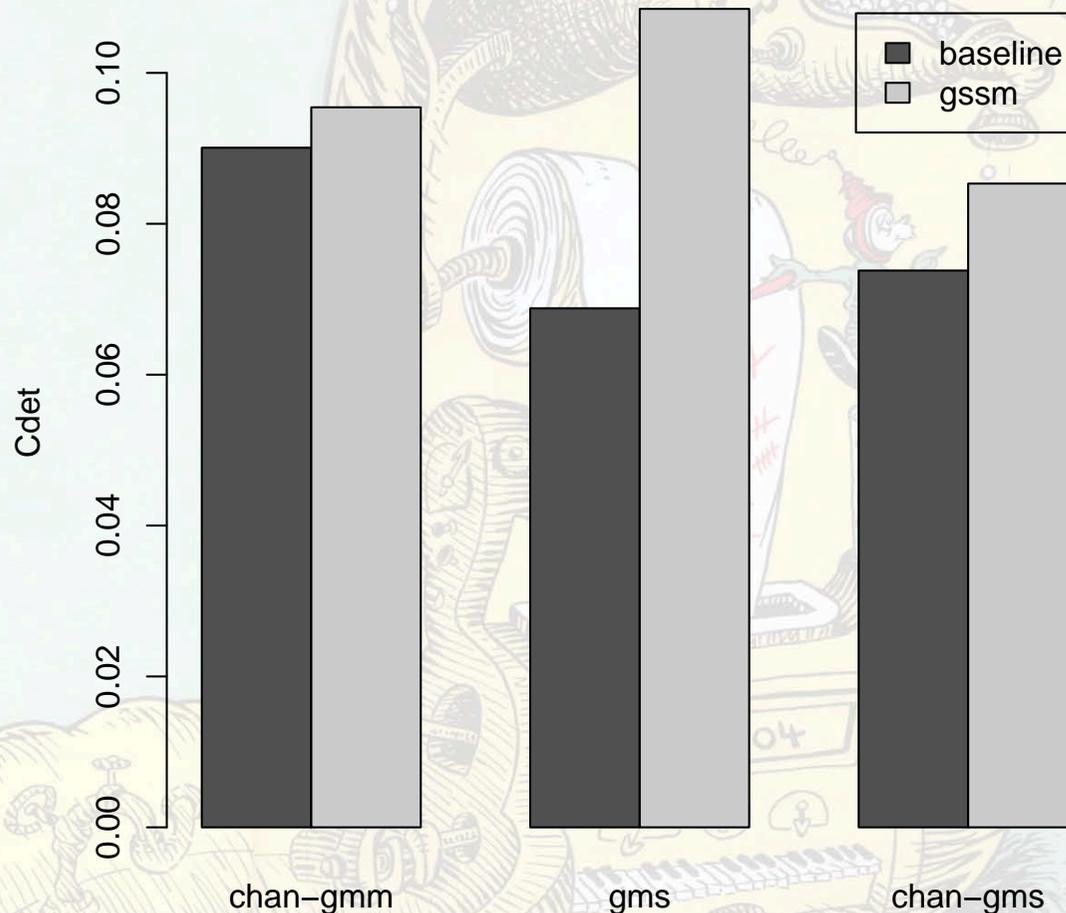
GMM means

+ −

$N_{ch} \times N_L$
Scores

SVM

# Results: Sparse Training Performance

SLD: Single Language Detector
GMS: GMM means SVM

- 30–60 hours per language for SLD

- 1.9–7.6 hours per language for LDA
  - collection of NIST trial sets '96–'03

- Observations

  - GMS best baseline

  - Chan-GMM most robust

  - Chan-GMS best sparse training

# Results: Effect of number of Single Language Detectors

- 'Columns' in LDA matrix

- random selection of $r$ columns per language
  - $r = 1\ldots6$
    - average 10 runs

- Chan-GMS system

- sparse training constant hit

**Effect of number of SLDs**

# Results: effect of sparse training size

- 'Rows' in LDA matrix

- Fraction of LDA training trials retained
  - $2^{-5} \ldots 2^0$
  - random selection
  - average 10 runs

- GMS
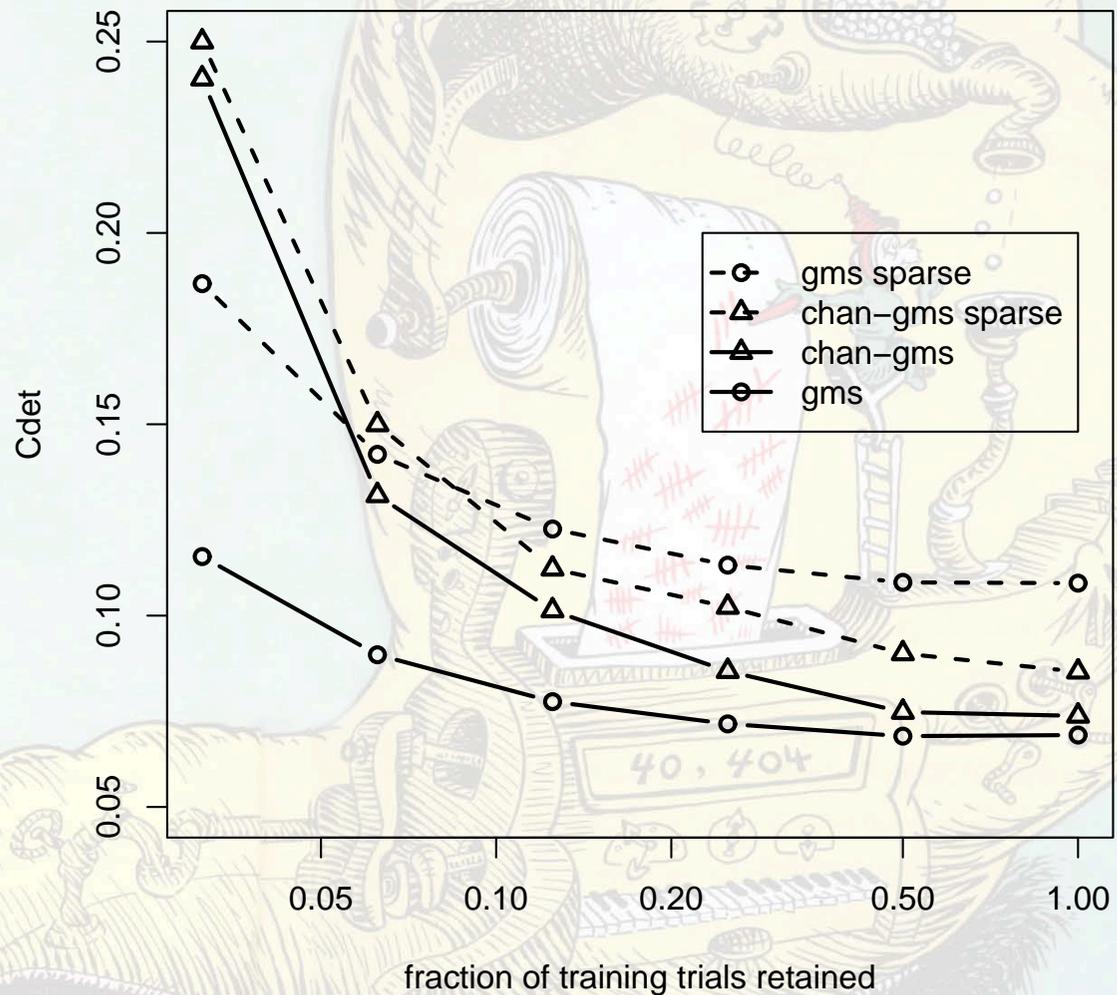  - large hit sparse
  - less hit by training size

- Chan-GMS
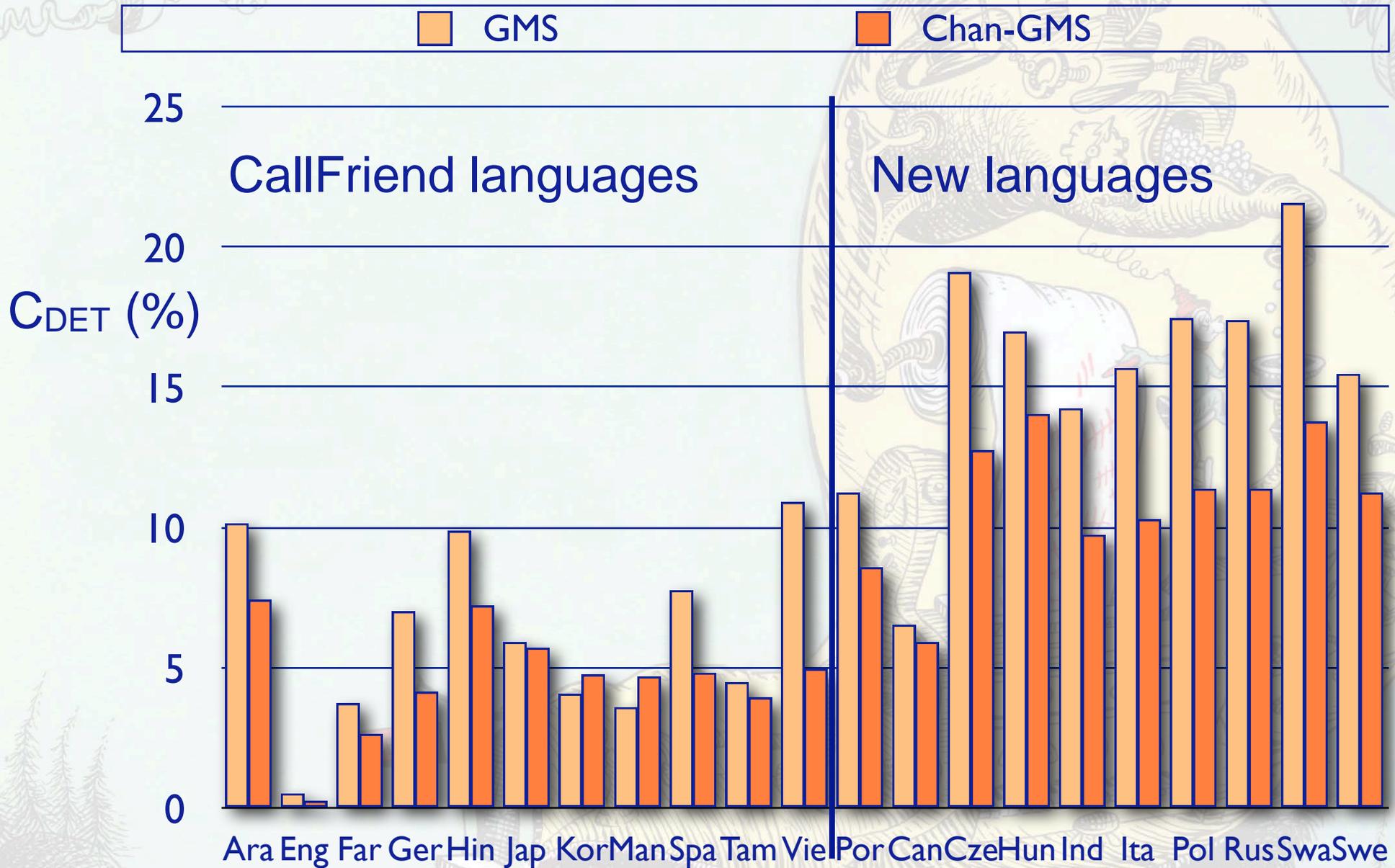  - smaller hit sparse
  - more hit by LDA

**Effect of LDA training size**

# Final test: independent data collection

- Use CSLU22 data collection as independent test
  - 21 languages
  - 2000+ speakers
  - Superset of LRE-2005 languages
  - 'story' sentences, 37s mean duration
  - 10-fold cross validation LDA-train / test

- ~ 54 min LDA training / language

- Full CallFriend training for SLDs

# Results CSLU22 per language



**Legend:** GMS (light orange), Chan-GMS (orange)

$C_{DET}$ (%)

CallFriend languages | New languages

Languages (x-axis): Ara Eng Far Ger Hin Jap Kor Man Spa Tam Vie Por Can Cze Hun Ind Ita Pol Rus Swa Swe

# Results CSLU: in/out set SLDs



Legend: ■ mean CallFriend  ■ mean New  ■ LRE-05 languages

$C_{DET}$ (%)

GMS: mean CallFriend 6.2, mean New 15.1, LRE-05 languages 5.2

Chan-GMS: mean CallFriend 4.6, mean New 10.6, LRE-05 languages 4.5

# Conclusions

- LDA can model new language for LID quite efficiently
    - very fast training of LDA
    - ~ 1 hour of training data gives $C_{DET}$ within factor ~ 2

- Generative GMMs seems more robust for missing SLD
    - but baseline performance is worse than discriminative GMS
    - rely more on back-end, anyway

- More SLDs in LDA
    - make LDA more robust for new language missing in SLDs
    - need more training data for LDA
        - *including* new language

- Discriminative channel-dependent GMS trade-off between
    - good baseline performance
    - fair robustness for language missing from SLDs