

Classifier Subset Selection and Fusion for Speaker Verification

**Filip Sedlák¹, Tomi Kinnunen¹,
Ville Hautamäki², Kong-Aik Lee², Haizhou Li^{1,2}**

¹University of Eastern Finland, School of
Computing, Joensuu, Finland

²Human Language Technology Department,
Institute of Infocomm Research (I²R), A*STAR,
Singapore

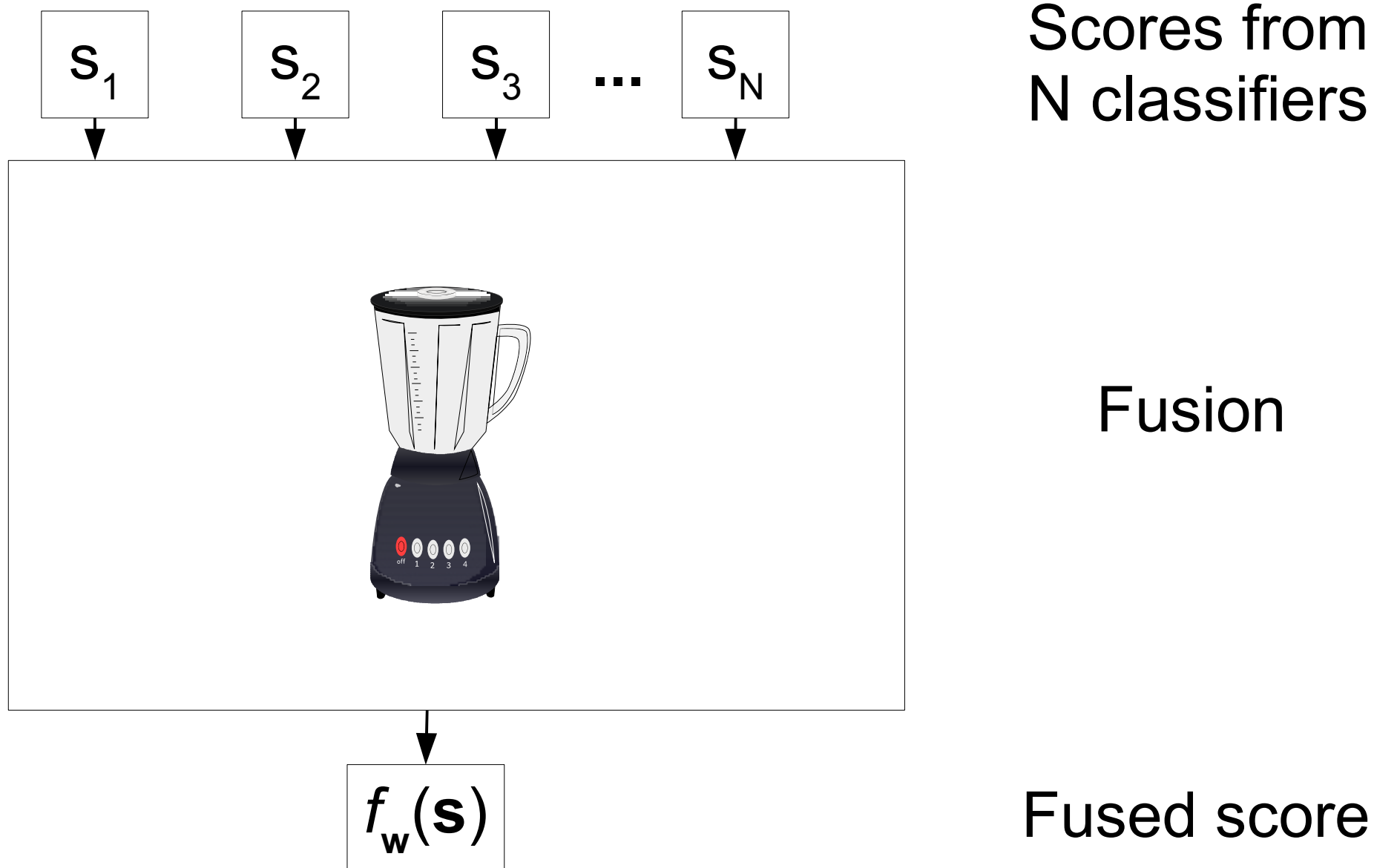
Classifier Fusion



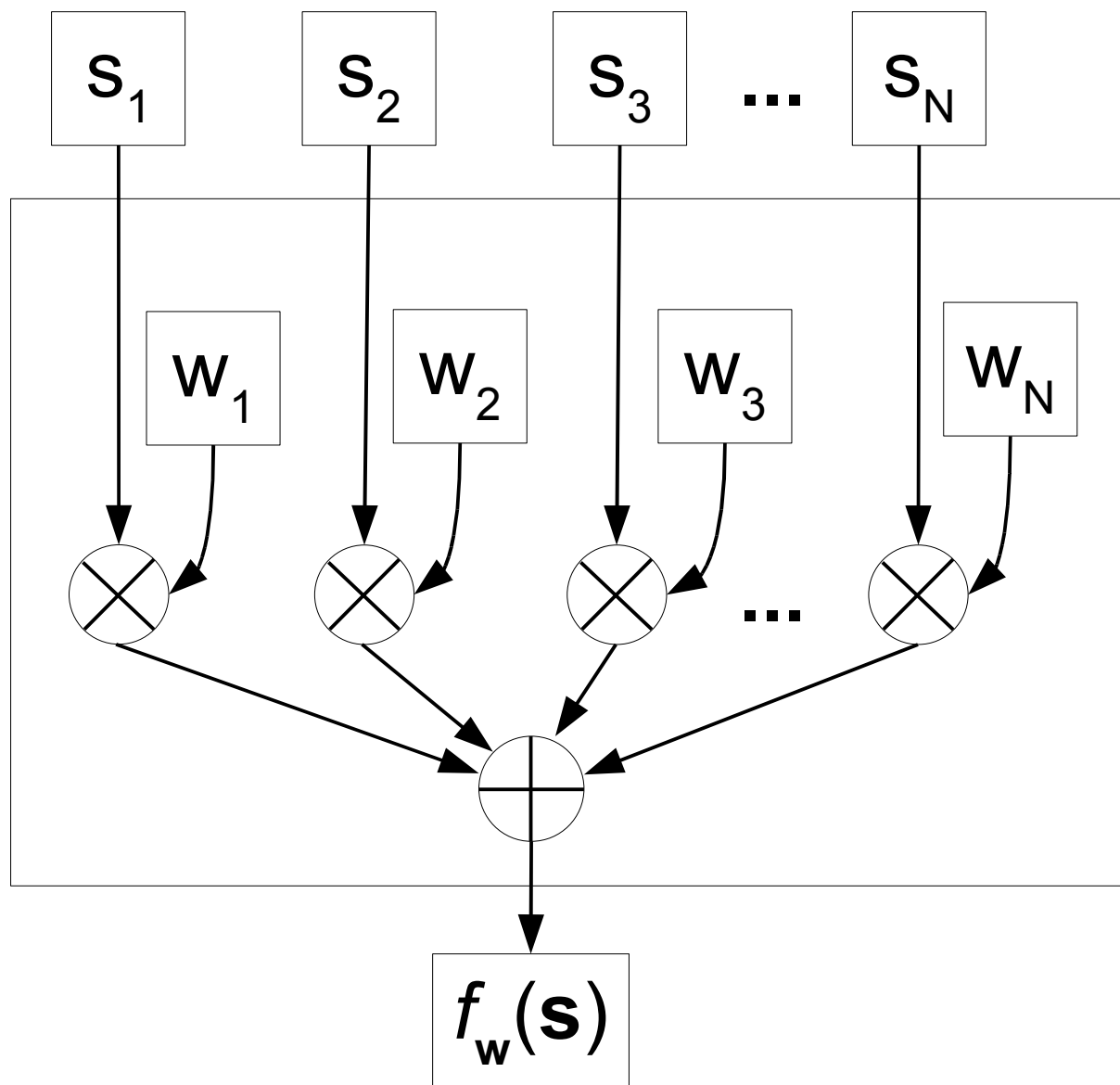
Verdict?



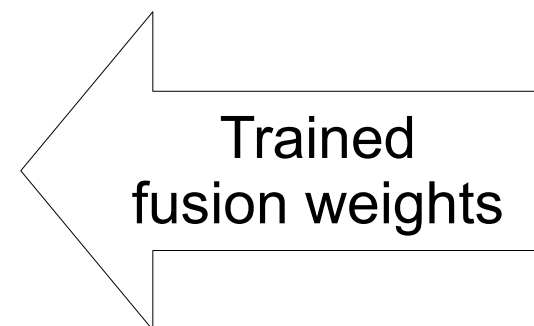
Classifier Fusion



Linear Fusion



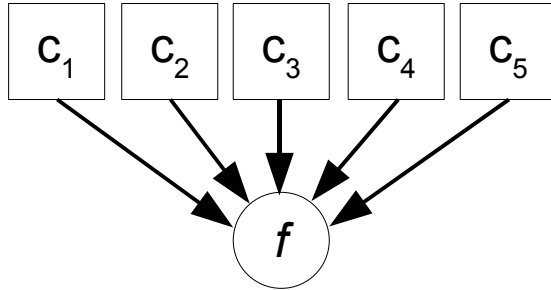
Scores from
N classifiers



$$f_w(\mathbf{s}) = \sum_{k=1}^N w_k s_k$$

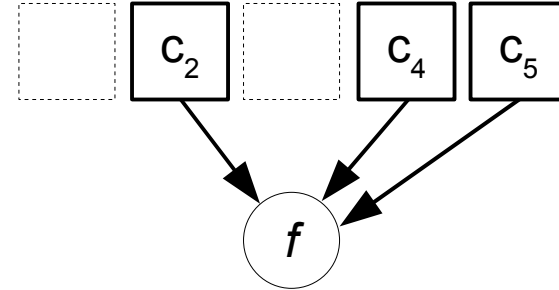
Full Set vs. Subset

Full set of $N=5$ classifiers



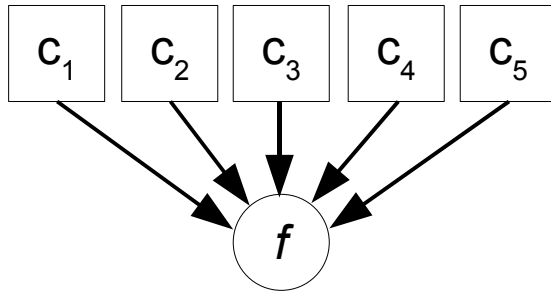
vs.

Subset of size 3



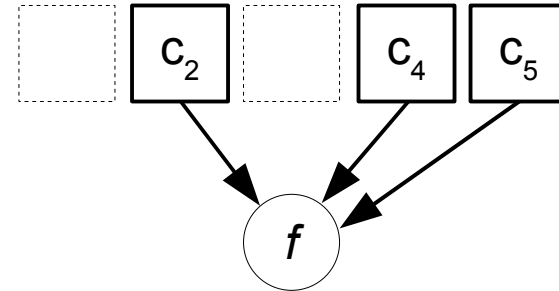
Full Set vs. Subset

Full set of $N=5$ classifiers



vs.

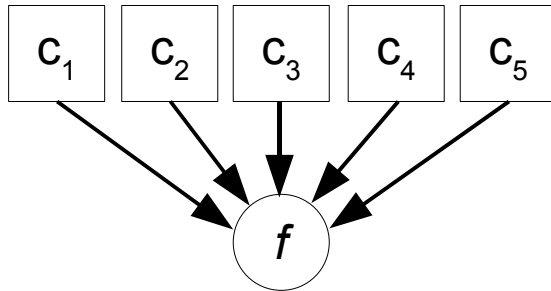
Subset of size 3



- + straightforward
- + computationally efficient
- possibly over-fitting if N is large

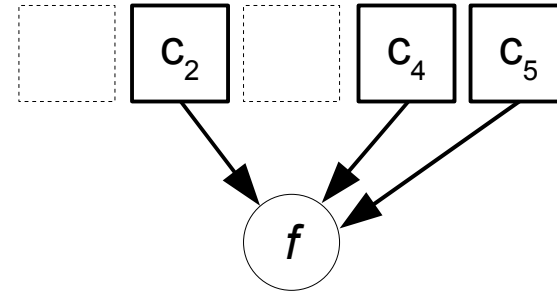
Full Set vs. Subset

Full set of $N=5$ classifiers



vs.

Subset of size 3

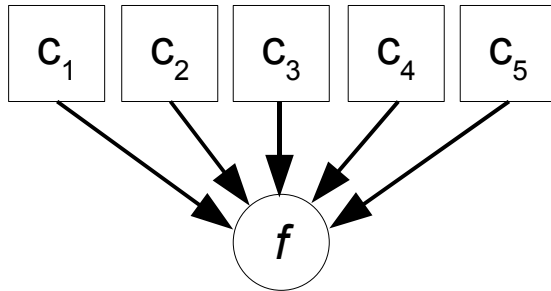


- + straightforward
- + computationally efficient
- possibly over-fitting if N is large

- + possibly better generalization
- requires subset selection

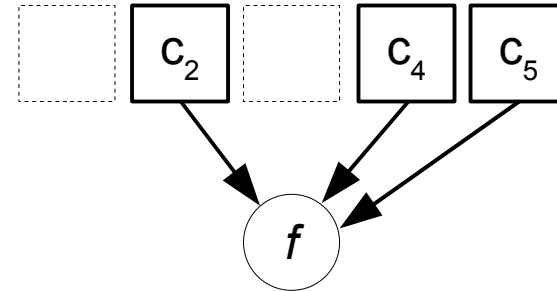
Full Set vs. Subset

Full set of $N=5$ classifiers



vs.

Subset of size 3

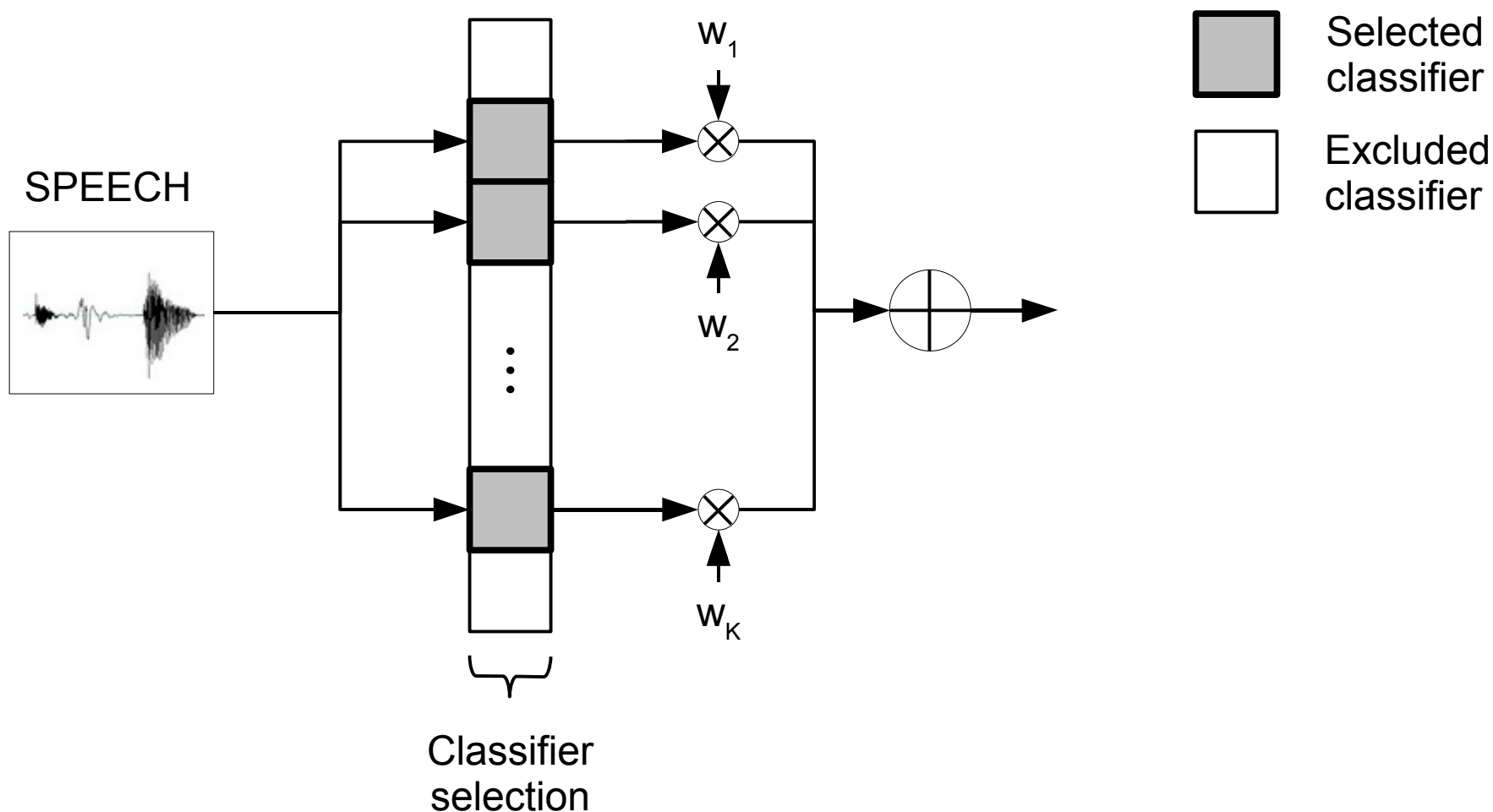


- + straightforward
- + computationally efficient
- possibly over-fitting if N is large

- + possibly better generalization
- requires subset selection

Can a subset fusion give better performance than the full-set?

Joint Classifier Selection and Fusion



Joint Classifier Selection and Fusion

Train **S-Cal** [1]
score warping
for each
individual
subsystem

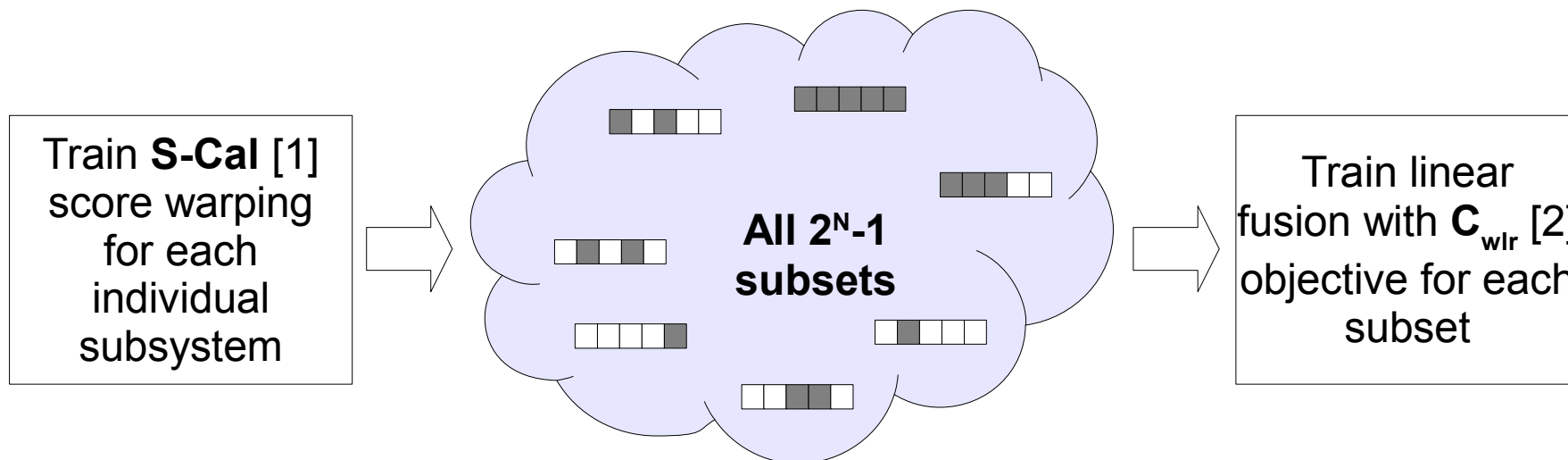
$$\text{SCal}(s) = \log \frac{(\text{logit}^{-1} \alpha)(e^{as+b} - 1) + 1}{(\text{logit}^{-1} \beta)(e^{as+b} - 1) + 1}$$

$$\text{logit}(x) = \log \frac{x}{1-x}$$

$$C_{\text{llr}} = \frac{1}{K} \sum_{i=1}^K \log(1 + e^{-S_{\text{Target}}}) + \frac{1}{L} \sum_{j=1}^L \log(1 + e^{S_{\text{Non-Target}}})$$

[1] Niko Brummer and Johan du Preez, "Application-Independent Evaluation of Speaker Detection", Computer Speech and Language, 2005.

Joint Classifier Selection and Fusion

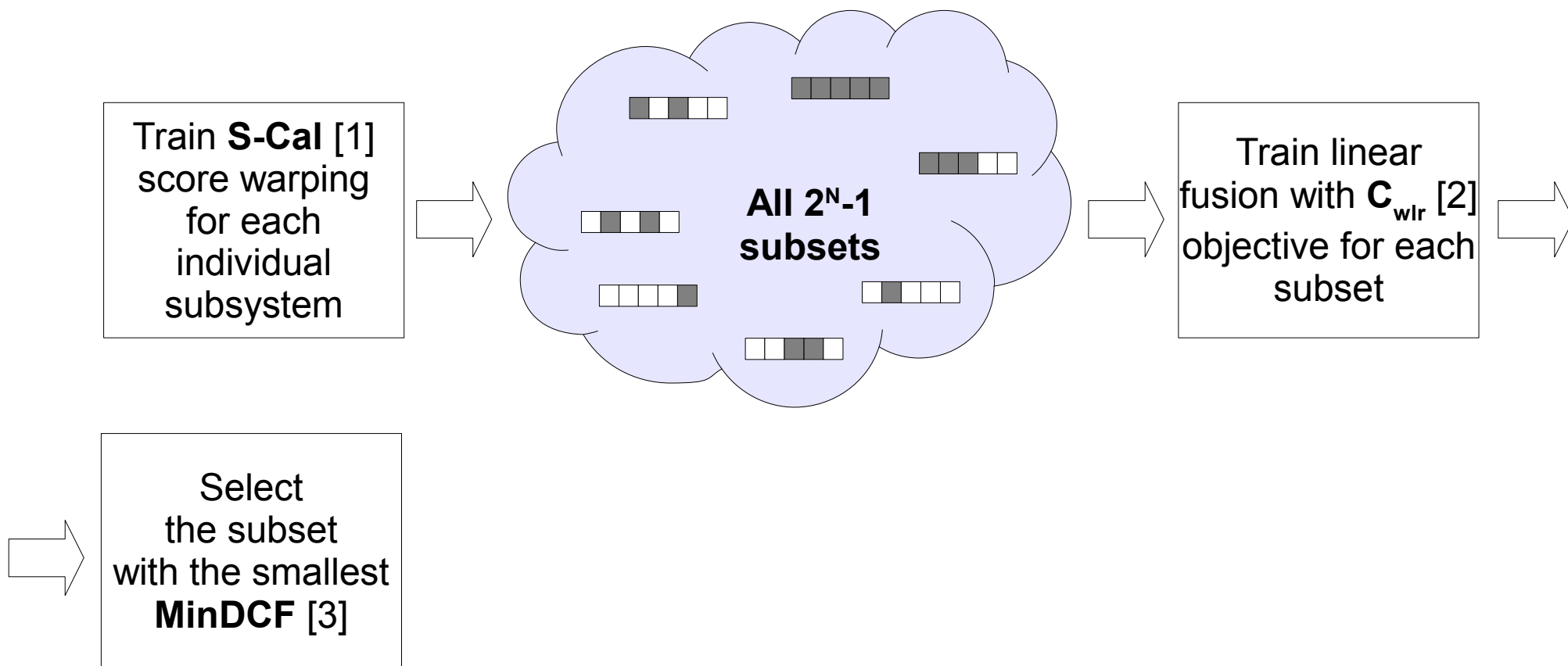


$$C_{wlr} = \frac{P}{K} \sum_{i=1}^K \log(1 + e^{-S_{Target} - \text{logit}(P)}) + \frac{1-P}{L} \sum_{j=1}^L \log(1 + e^{S_{Non-Target} + \text{logit}(P)})$$

$$P = \text{logit}^{-1} \left(\text{logit}(P_{Target}) + \log \frac{C_{Miss}}{C_{FA}} \right)$$

Here $C_{Miss} = 1$, $C_{FA} = 1$, $P_{Target} = 0.001$ (i.e. the 'new' NIST cost function).

Joint Classifier Selection and Fusion

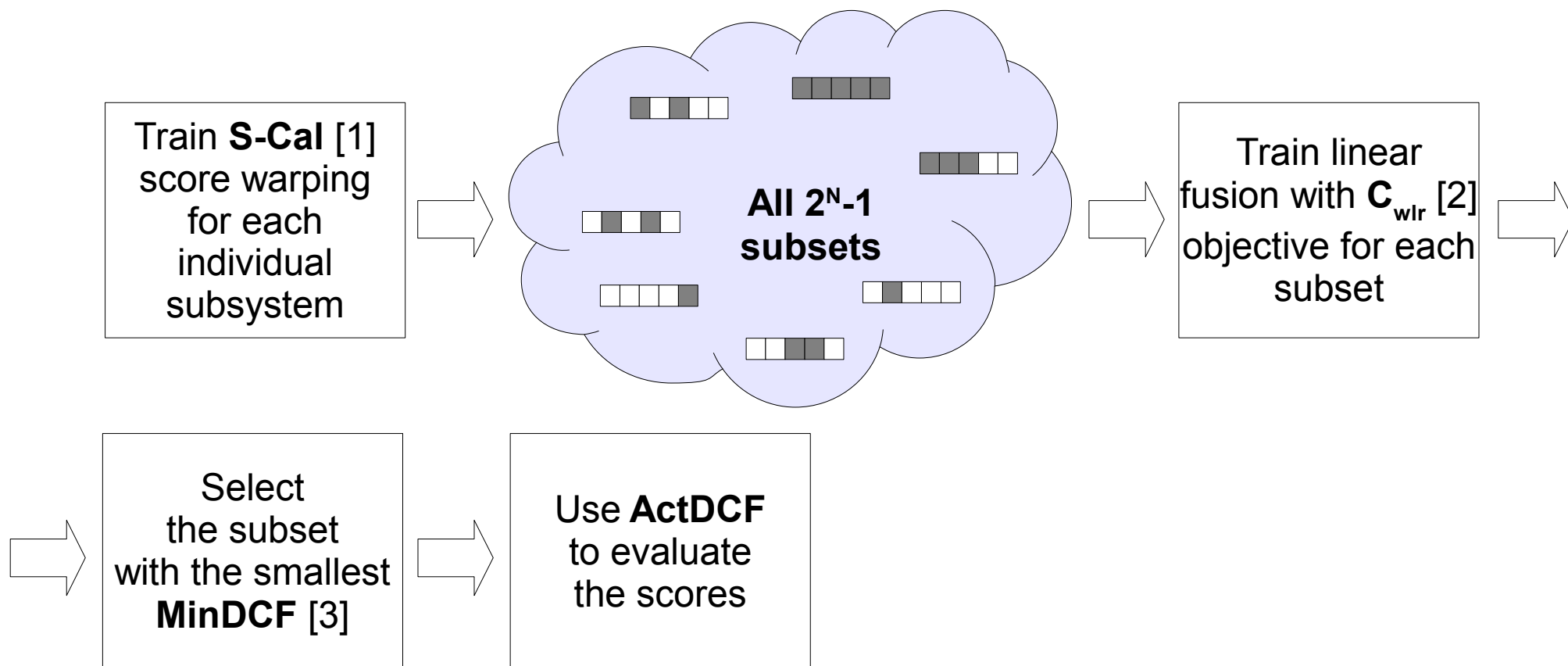


$$\text{DCF}(\theta) = C_{\text{Miss}} P_{\text{Miss}}(\theta) P_{\text{Target}} + C_{\text{FA}} P_{\text{FA}}(\theta) (1 - P_{\text{Target}})$$

$$\text{MinDCF} = \min_{\theta} \text{DCF}(\theta)$$

Here $C_{\text{Miss}} = 1$, $C_{\text{FA}} = 1$, $P_{\text{Target}} = 0.001$ (i.e. the 'new' NIST cost function).

Joint Classifier Selection and Fusion



$$DCF(\theta) = C_{Miss} P_{Miss}(\theta) P_{Target} + C_{FA} P_{FA}(\theta) (1 - P_{Target})$$

$$ActDCF = DCF(\theta_{Trainset})$$

Here $C_{Miss} = 1$, $C_{FA} = 1$, $P_{Target} = 0.001$ (i.e. the 'new' NIST cost function).

Base Classifiers (I4U, NIST 2010)

	Classifier	Feature	Evalset1 EER (%)	Evalset2 EER (%)
1	GMM-UBM	PLP	3.95	4.99
2	Joint Factor Analysis	PLP	4.24	4.12
3		PLP	4.24	3.75
4		LPCC	4.59	5.74
5	GMM-SVM	PLP	5.65	5.49
6	Kullback–Leibler Divergence	MFCC	4.99	4.37
7		LPCC	6.45	5.37
8		MLF	5.81	4.74
9		LPCC	4.24	6.52
10		SWLP	10.20	5.87
11	GMM-SVM Feature Transformation	PLP	8.13	6.12
12	GMM-SVM Bhattacharyya Distance	PLP	5.40	3.03

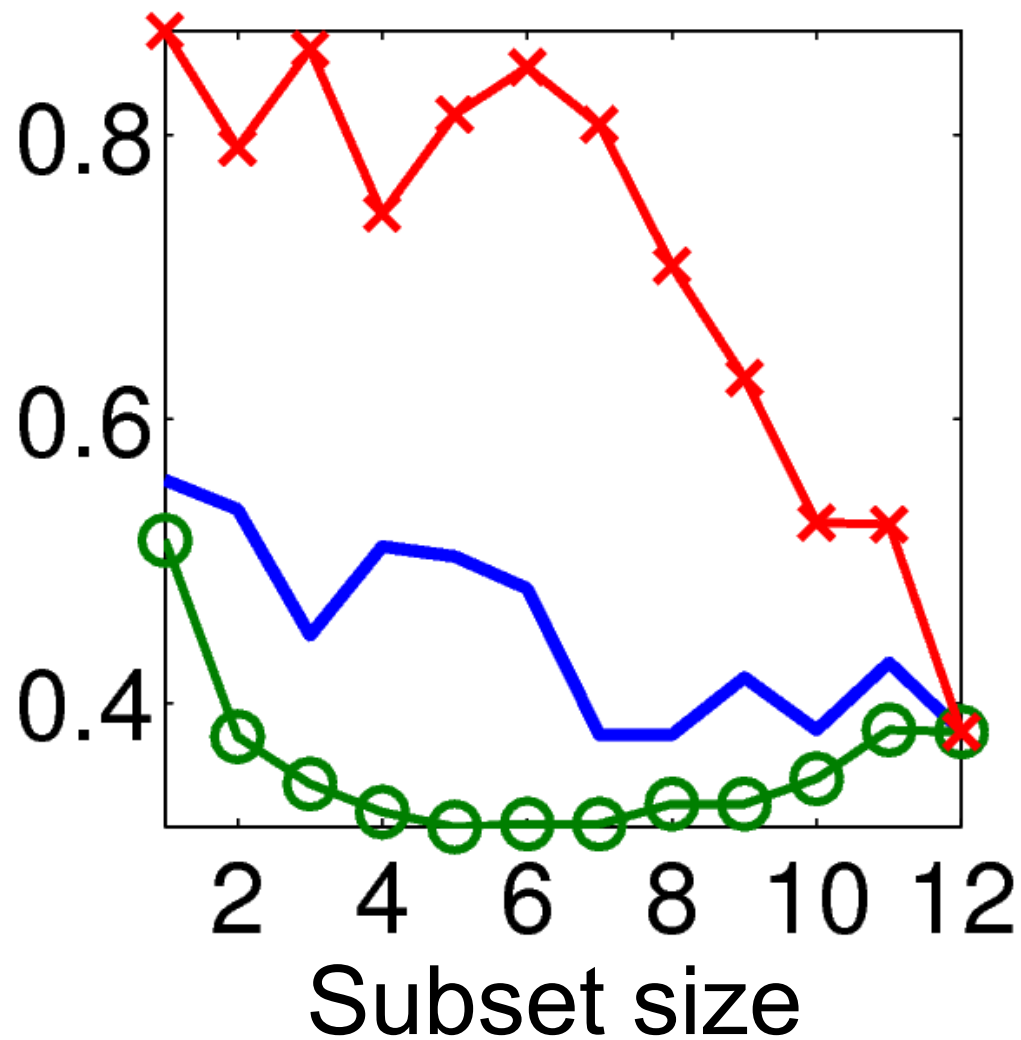
All datasets:
interview-telephone,
female trials

Trainset,
Evalset1:
NIST SRE 2008

Evalset2:
NIST SRE 2010

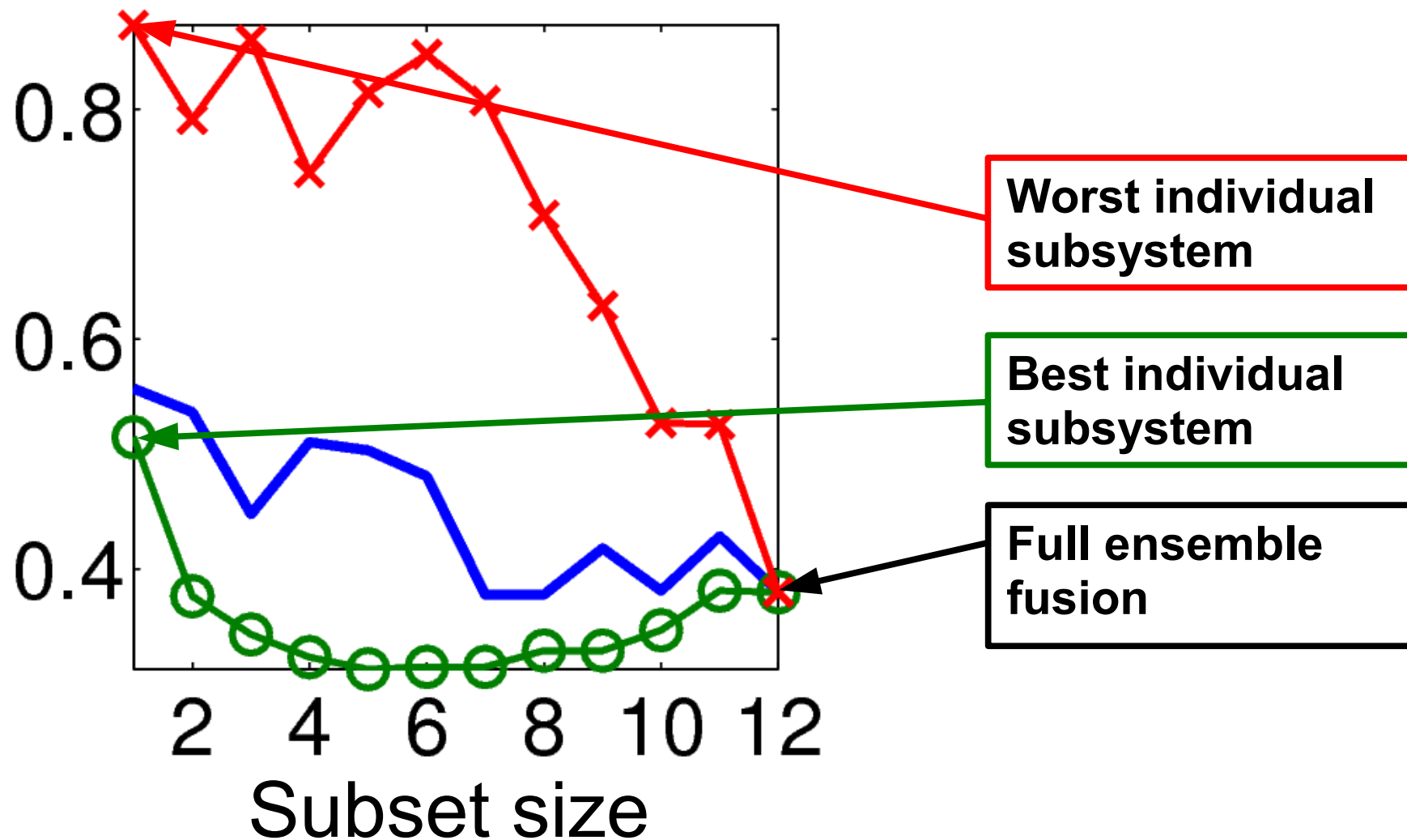
Error Bounds on the NIST2008 (Evalset1)

1000 x ActDCF



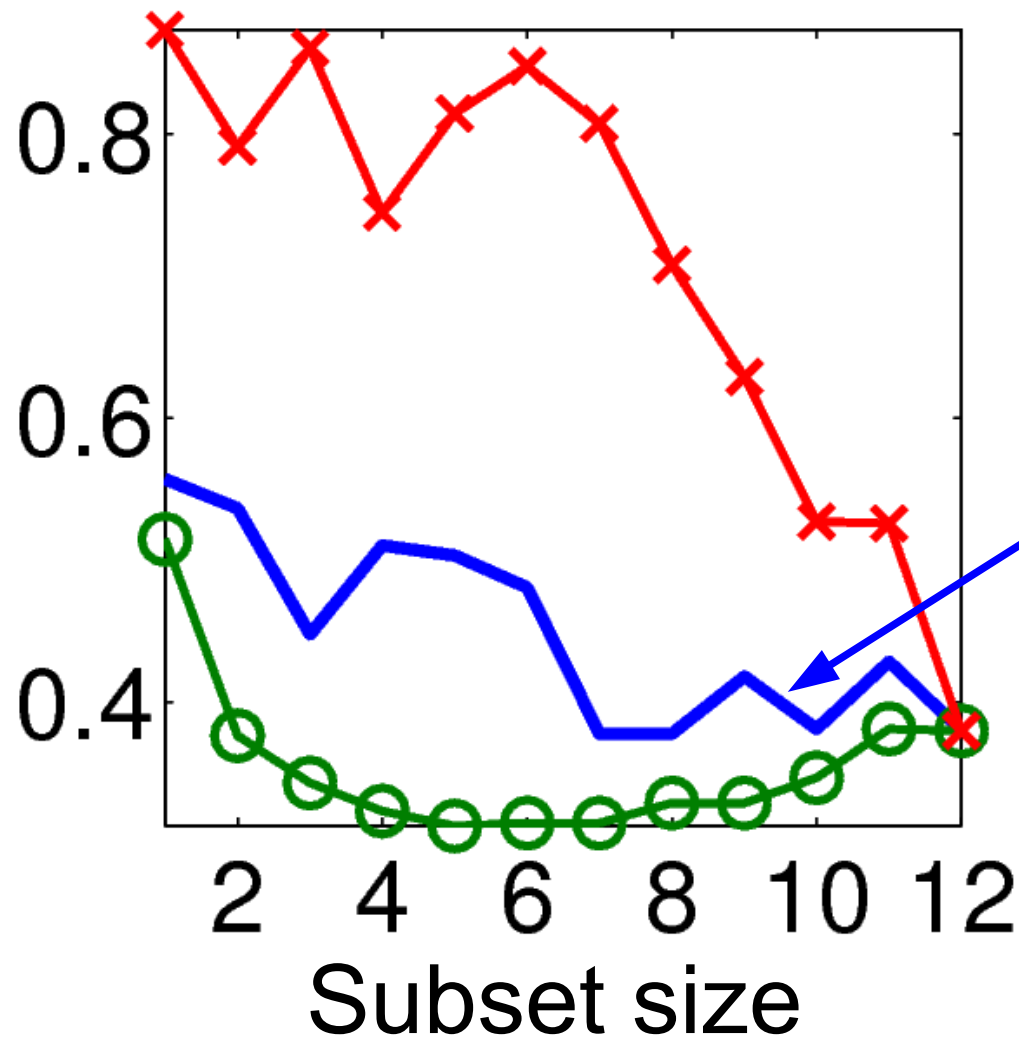
Error Bounds on the NIST2008 (Evalset1)

1000 x ActDCF



Error Bounds on the NIST2008 (Evalset1)

1000 x ActDCF

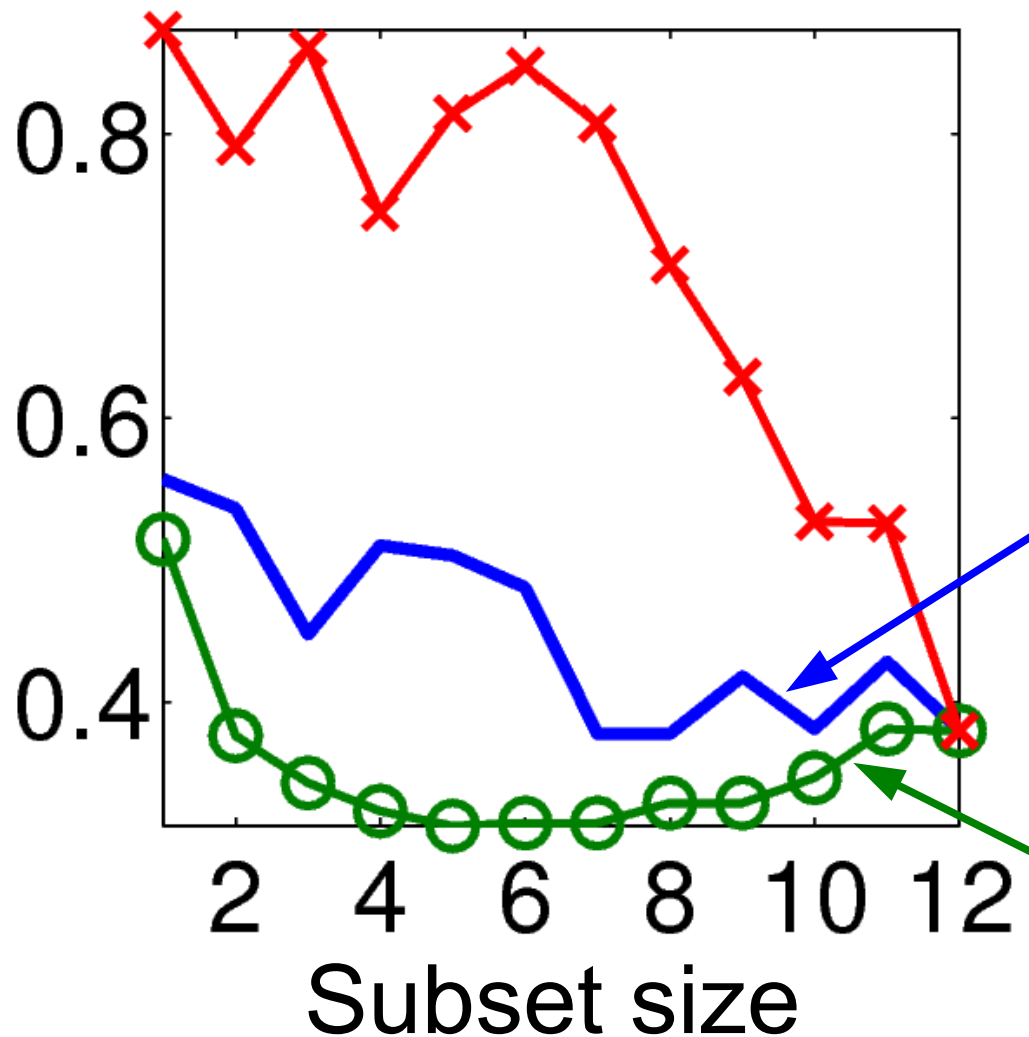


Realistic use-case

Fusion is trained on the **trainset**, subset is selected based on the **best performance on the trainset**.

Error Bounds on the NIST2008 (Evalset1)

1000 x ActDCF

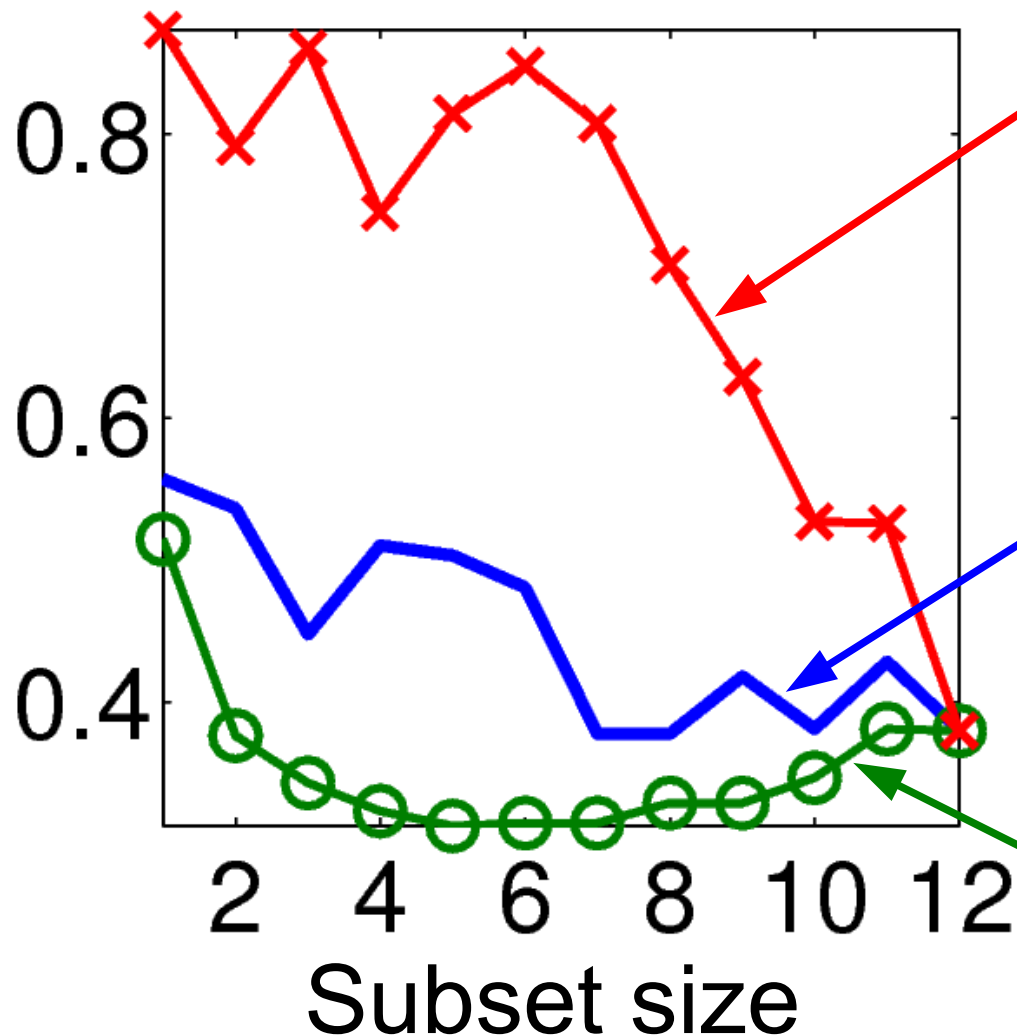


Realistic use-case
Fusion is trained on the **trainset**, subset is selected based on the **best performance on the trainset**.

Best subset selection oracle
Fusion is trained on the **trainset**, subset is selected based on the **best performance on the evalset1**.

Error Bounds on the NIST2008 (Evalset1)

1000 x ActDCF



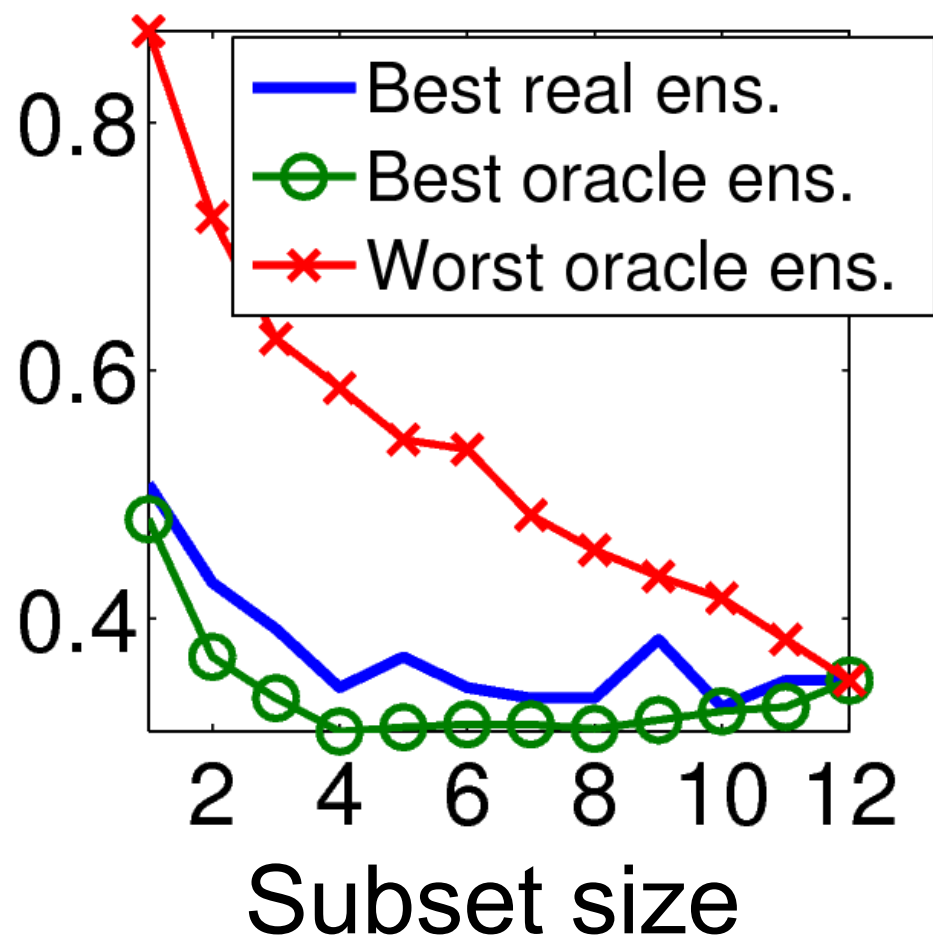
Worst subset selection oracle
Fusion is trained on the **trainset**, subset is selected based on the **worst** performance on the **evalset1**.

Realistic use-case
Fusion is trained on the **trainset**, subset is selected based on the **best** performance on the **trainset**.

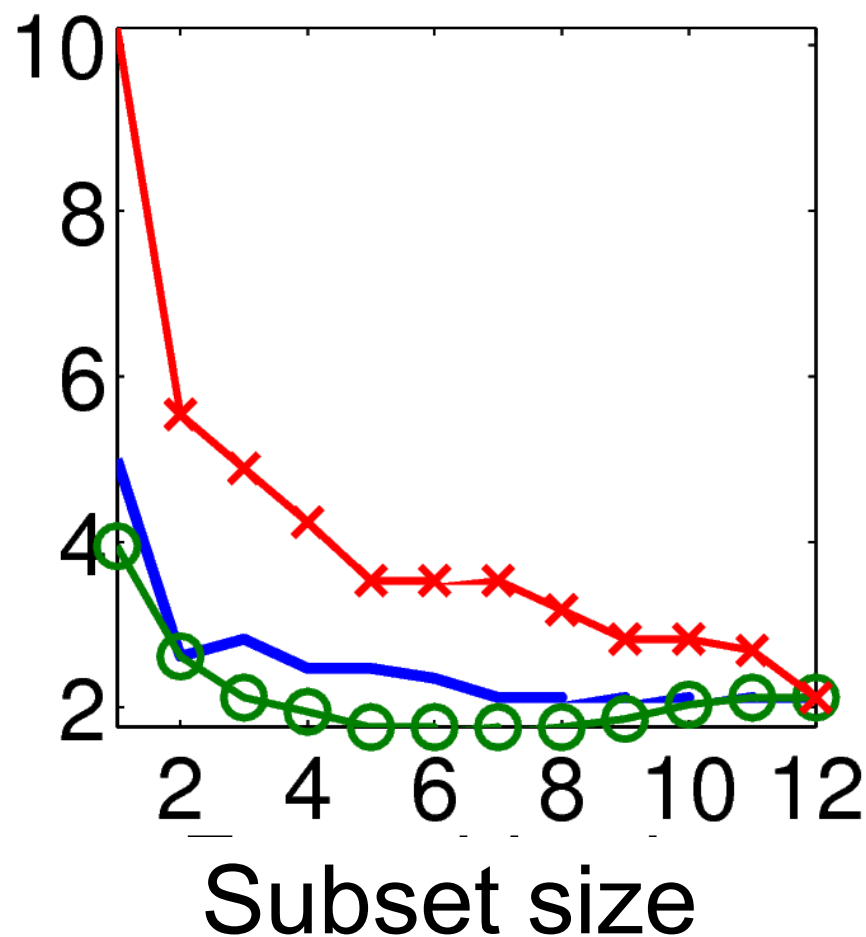
Best subset selection oracle
Fusion is trained on the **trainset**, subset is selected based on the **best** performance on the **evalset1**.

Error Bounds on the NIST2008 (Evalset1)

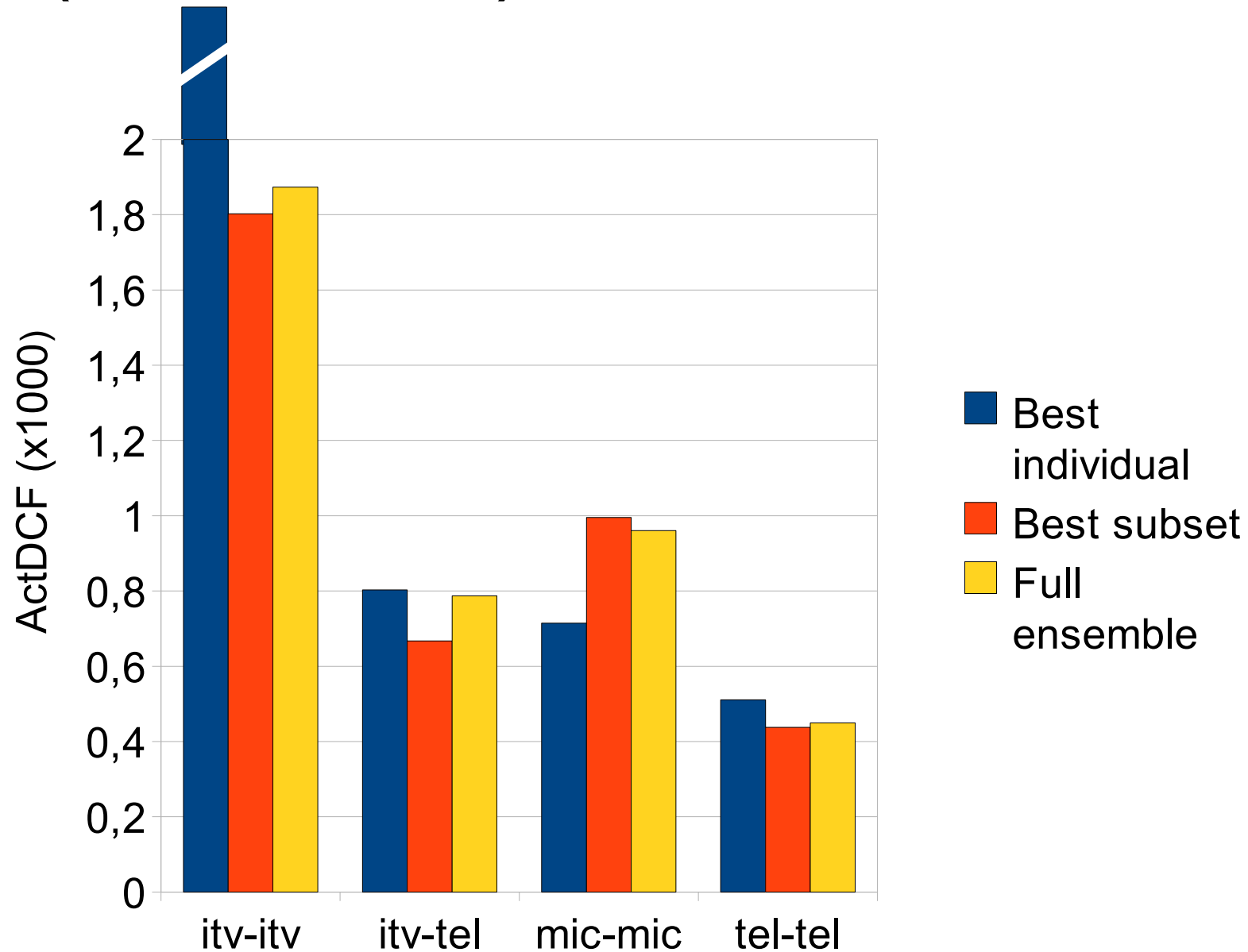
1000 x MinDCF



EER (%)

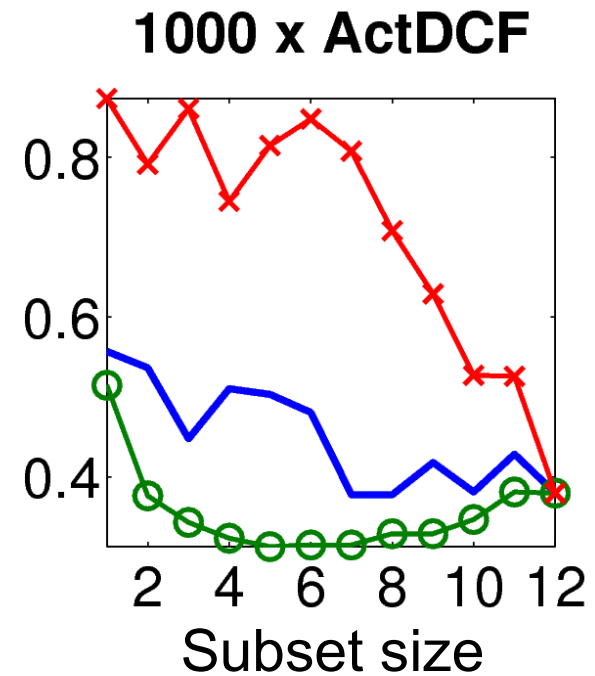
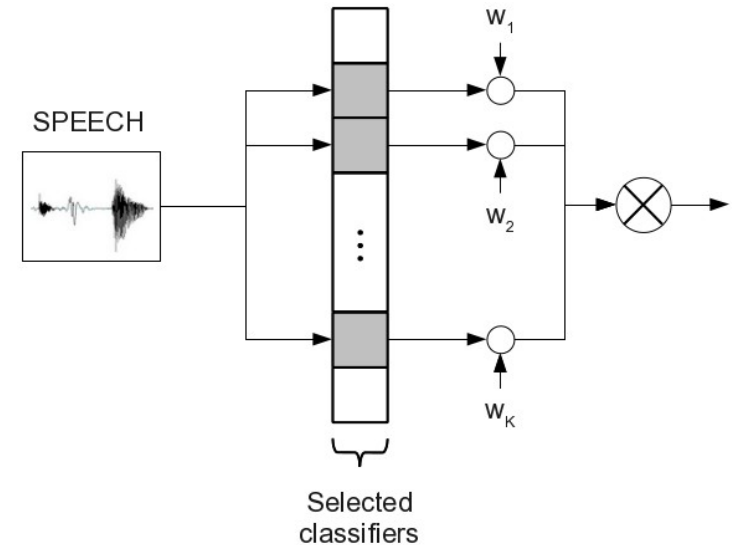


Subset Performance on Evalset2 (NIST2010), Pooled Genders



Conclusions

- Subset fusion has the potential to outperform the full set fusion.
- Further study should focus on subset selection methods.



Following are support slides...

Score Sets

	Trainset	Devset	Evalset
Trials	NIST2008 itv-tel female	NIST2008 itv-tel female	NIST2010 itv-tel female
Target	263	283	801
Non-target	27 315	27 195	30 254

Error Bounds

- Best individual base system
- Full set fusion

Weights trained on	Best subset selected on	Performance evaluated on	
Trainset	Trainset	Devset/Evalset	'Real'
Trainset	Devset/Evalset	Devset/Evalset	'Best Real'
Devset/Evalset	Devset/Evalset	Devset/Evalset	'Best Oracle'