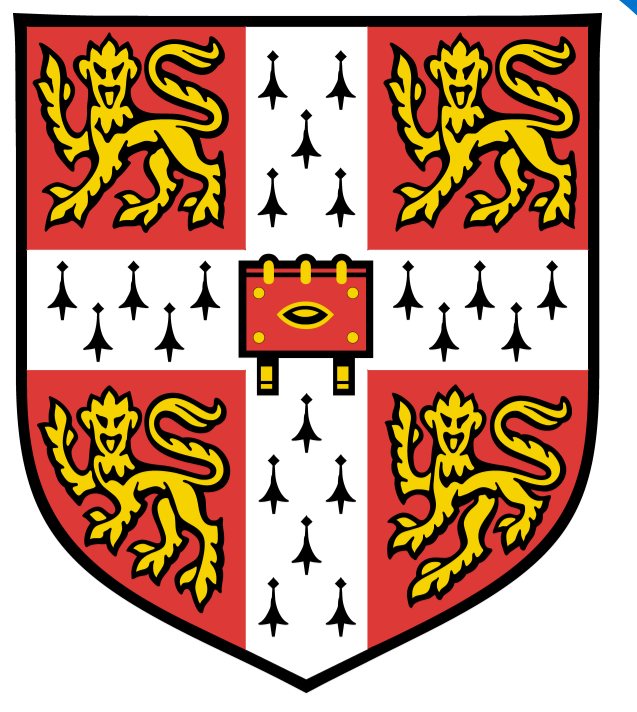


# ASR teacher-student training and ensemble target diversity

Jeremy H. M. Wong and Mark J. F. Gales

Department of Engineering, University of Cambridge

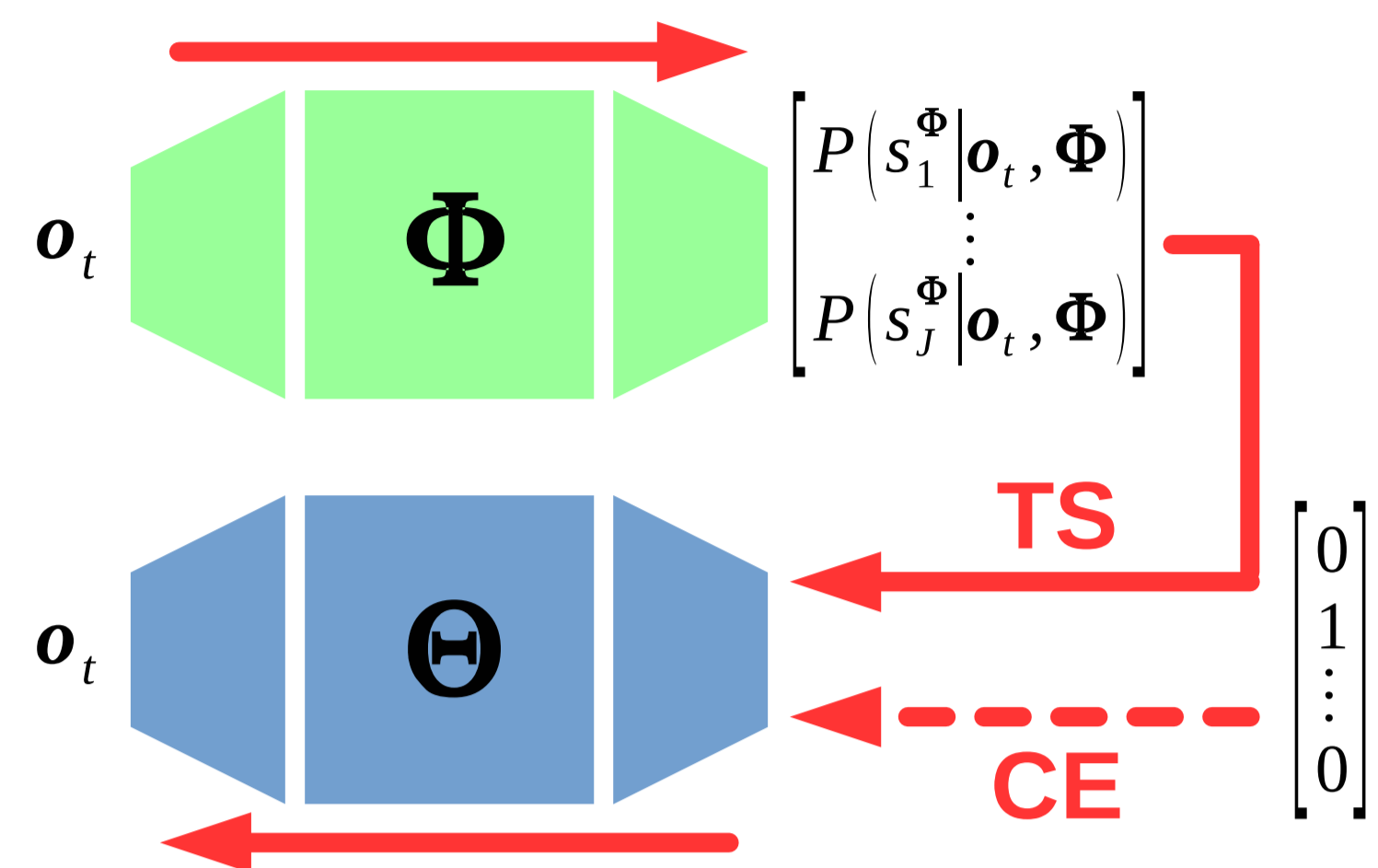
jhmw2@cam.ac.uk, mjfg@eng.cam.ac.uk



## MOTIVATION

- An ensemble with **target diversity** can give good combination gains.
- Improve recognition efficiency by training a student to emulate the ensemble.
- How to propagate information across different output targets?

## TEACHER-STUDENT TRAINING



- Train a single student model to emulate the ensemble behaviour.

- **Standard CE training:**

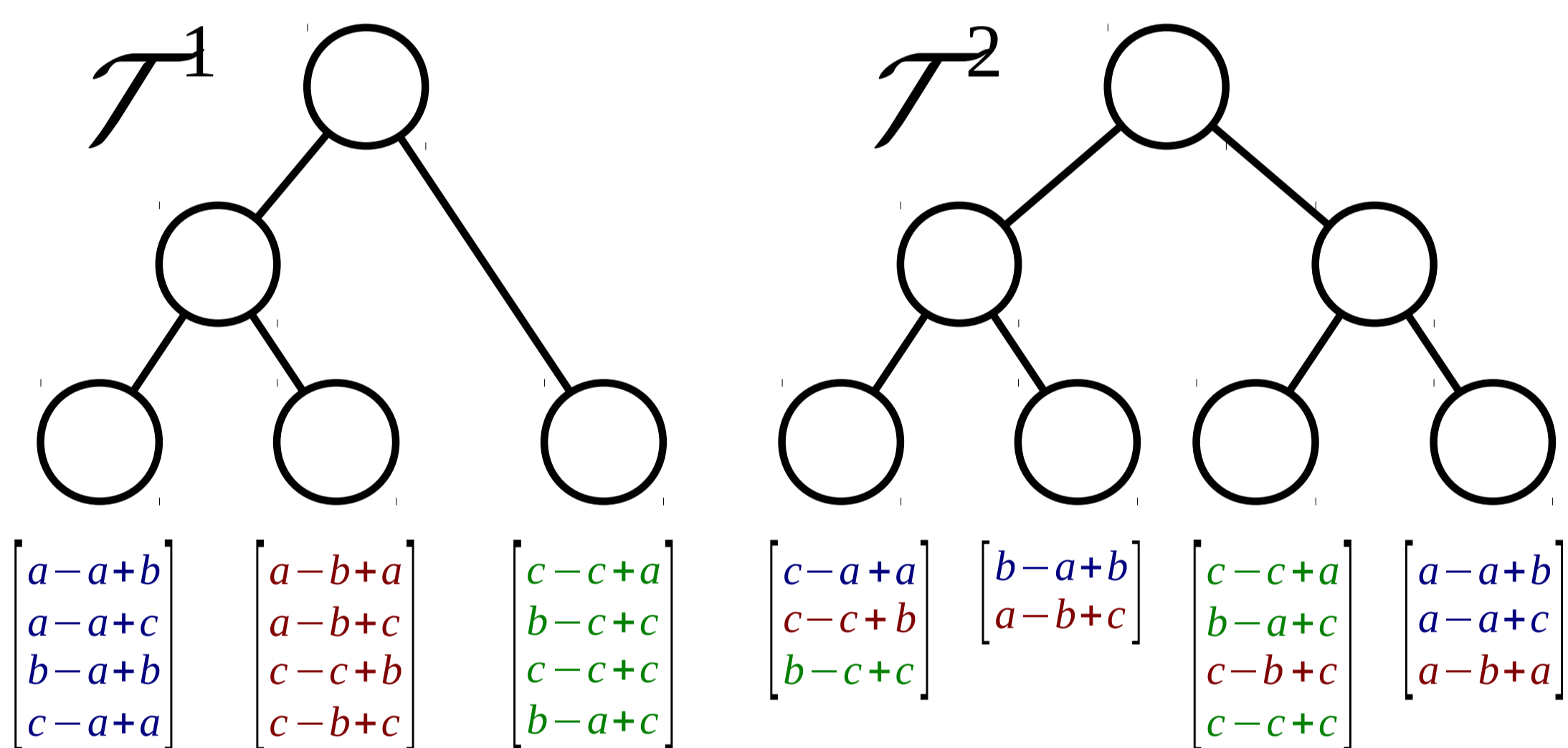
$$\mathcal{F}_{CE} = - \sum_t \sum_{s \in \mathcal{T}} \delta(s, s_t^*) \log P(s|o_t, \Theta)$$

- **Standard Teacher-Student (TS) training:**

$$\mathcal{F}_{TS} = - \sum_t \sum_{s \in \mathcal{T}} \sum_{m=1}^M \lambda_m P(s|o_t, \Phi^m) \log P(s|o_t, \Theta)$$

- Use only student model during recognition.
- Requires student's and teacher's outputs to have the same interpretations.

## OUTPUT TARGET DIVERSITY



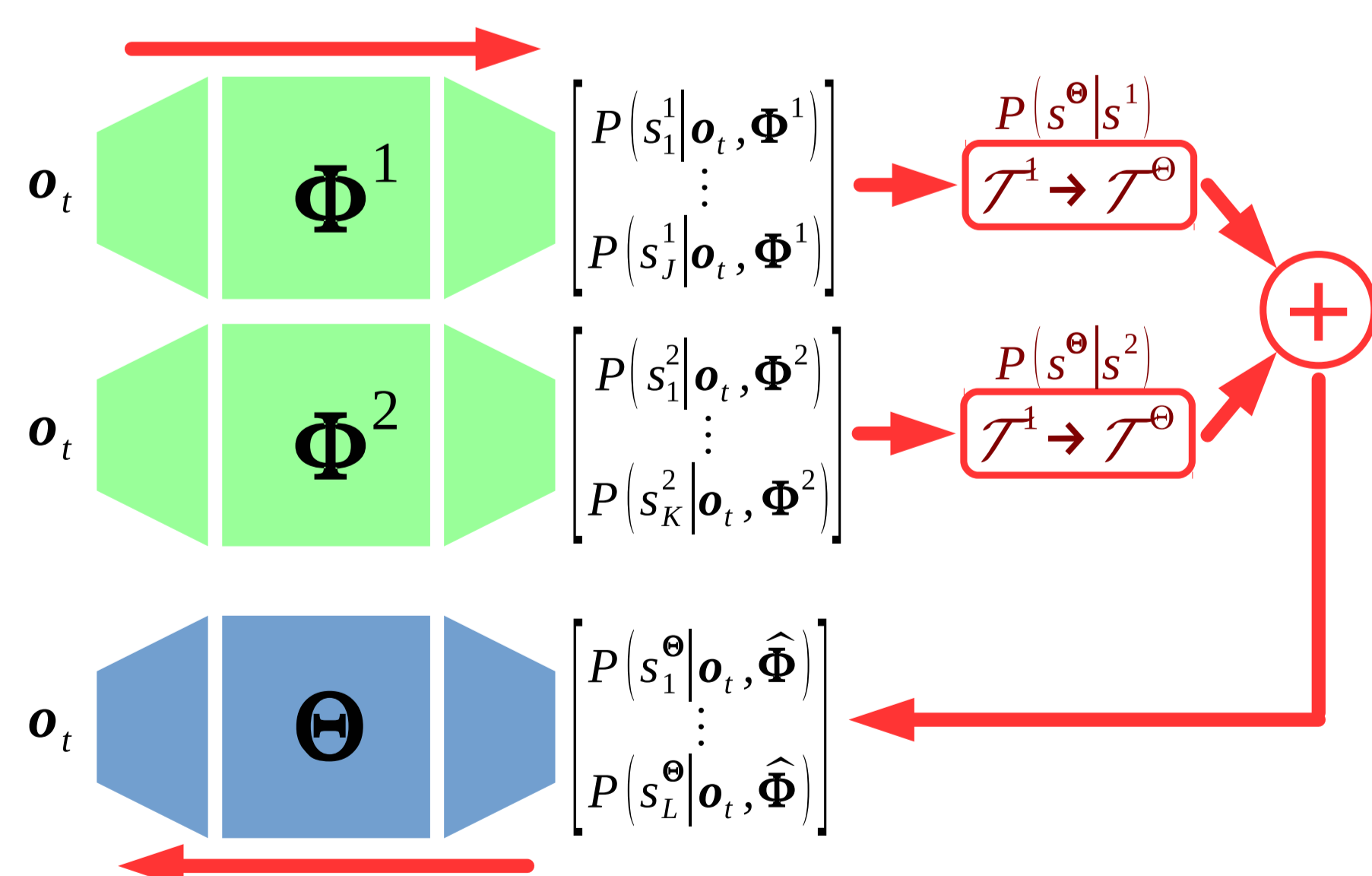
- Output targets are defined by a Phonetic Decision Tree (PDT)

$$s_c = \mathcal{T}(c)$$

- Generate ensemble by using a different PDT for each model.
- Models learn to discriminate between different sets of state clusters.
- Computational cost of ensemble combination:

	NN forward	lattice decode
hypothesis combine	$M$	$M$
frame combine	$M$	1
teacher-student	1	1

## POSTERIOR MAPPING



- When PDTs differ, train student by minimising logical context KL-divergence:

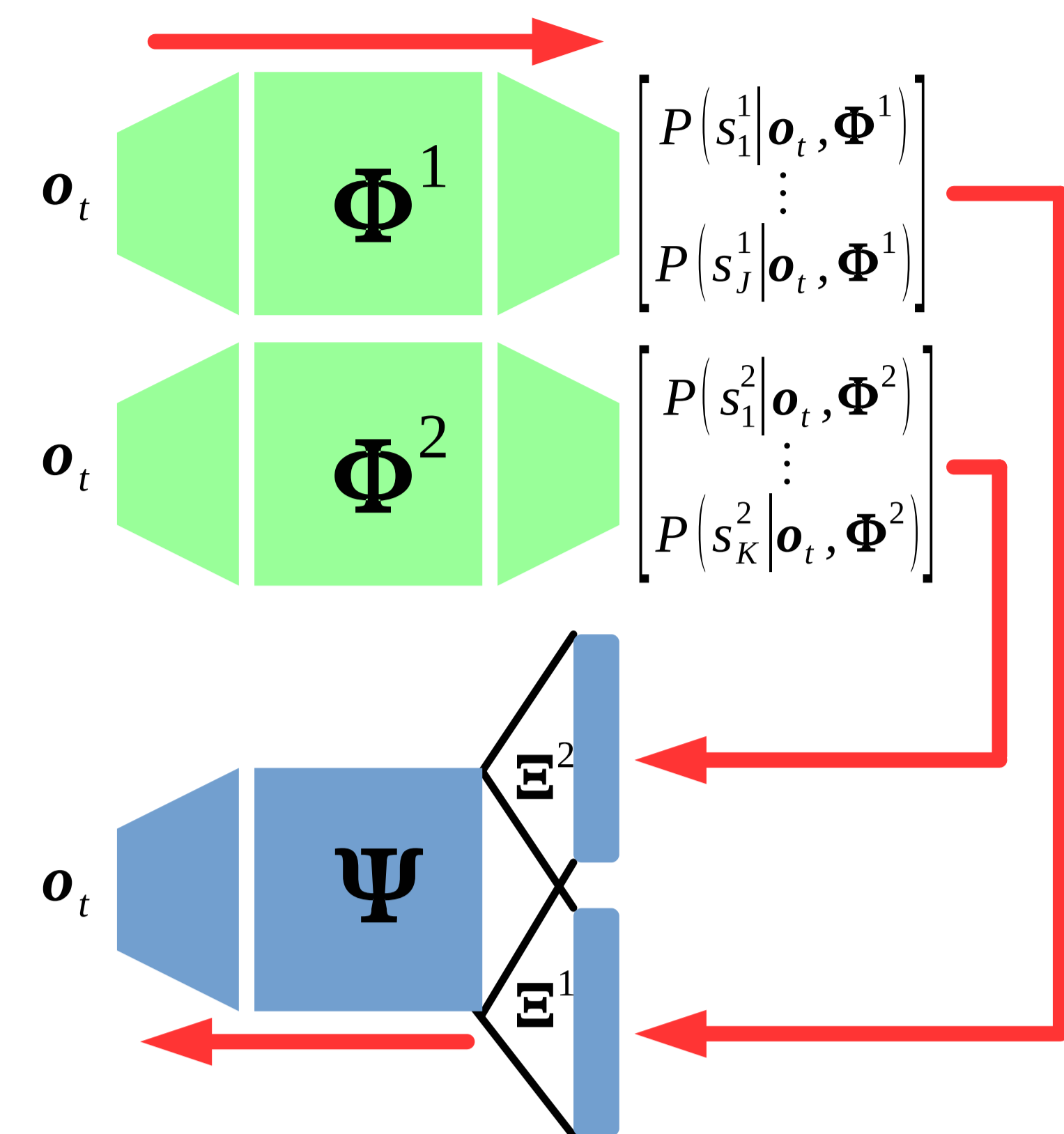
$$\mathcal{F}_{RF-TS} = - \sum_t \sum_{c \in \mathcal{C}} \sum_{m=1}^M \lambda_m P(c|o_t, \Phi^m) \log P(c|o_t, \Theta)$$

- Under mild assumptions, the criterion reduces to:

$$\mathcal{F}_{RF-TS} = - \sum_t \sum_{s^0 \in \mathcal{T}^0} \sum_{m=1}^M \lambda_m \sum_{s^m \in \mathcal{T}^m} P(s^0|s^m) P(s^0|o_t, \Phi^m) \log P(s^0|o_t, \Theta)$$

- $P(s^0|s^m)$  maps posteriors between PDTs.
- Can estimate  $P(s^0|s^m)$  from forced alignments.
- Student PDT size can be chosen independently of teacher PDTs.

## MULTI-TASK ARCHITECTURE



- Avoid mapping by using Multi-Task (MT) student.

- **Multi-task CE training:**

$$\mathcal{F}_{MT} = - \sum_t \sum_{m=1}^M \sum_{s^m \in \mathcal{T}^m} \delta(s^m, s_t^{m*}) \log P(s^m|o_t, \Psi, \hat{\Xi})$$

- **Multi-task teacher-student training:**

$$\mathcal{F}_{MT-TS} = - \sum_t \sum_{m=1}^M \sum_{s^m \in \mathcal{T}^m} P(s^m|o_t, \Phi^m) \log P(s^m|o_t, \Psi, \hat{\Xi})$$

## EXPERIMENTS

- **Datasets:**

- **207V: IARPA Babel Tok Pisin**  
\* 3 hours VLLP training set, 1000 PDT states
- **AMI: Augmented multi-party interaction**  
\* 81 hours IHM training set, 4000 PDT states
- **HUB4: English broadcast news**  
\* 144 hours training set, 6000 PDT states

- **Ensemble size = 4**

- Student and teachers have the same architecture.

## SINGLE MODEL PERFORMANCE

Dataset	Single model WER (%)				cross-WER (%)
	mean	best	worst	std dev	
207V	48.3	48.0	48.4	0.17	28.4
AMI	26.0	25.9	26.2	0.13	15.2
HUB4	9.3	9.2	9.4	0.10	7.0

- Measure diversity using cross-WER:

$$\text{cross-WER} = \frac{1}{M(M-1)} \sum_{m=1}^M \sum_{n \neq m}^M \text{WER}(\mathcal{H}^m, \mathcal{H}^n)$$

## ENSEMBLE PERFORMANCE

Dataset	ensemble	Combined WER (%)		
		hypothesis	frame	student
207V	separate	45.8	46.0	46.6
	MT	47.7	47.8	47.3
AMI	separate	24.5	24.6	24.6
	MT-TS	24.3	24.4	24.6
HUB4	separate	8.7	8.7	9.0
	MT-TS	8.8	8.7	8.9

- Single-output student can learn from teachers with different PDTs.
- Multi-task student is able to match the ensemble performance.

## CONCLUSIONS

- Proposed teacher-student method when output targets differ.
- Proposed multi-task teacher-student method.

## REFERENCES

- [1] J. Wong and M. Gales, "Student-teacher training with diverse decision tree ensembles", *Interspeech*, Sep 2017
- [2] J. Wong and M. Gales, "Multi-task ensembles with teacher-student training", *ASRU*, Dec 2017