

# Uncertainty: Knowing What You Don't Know

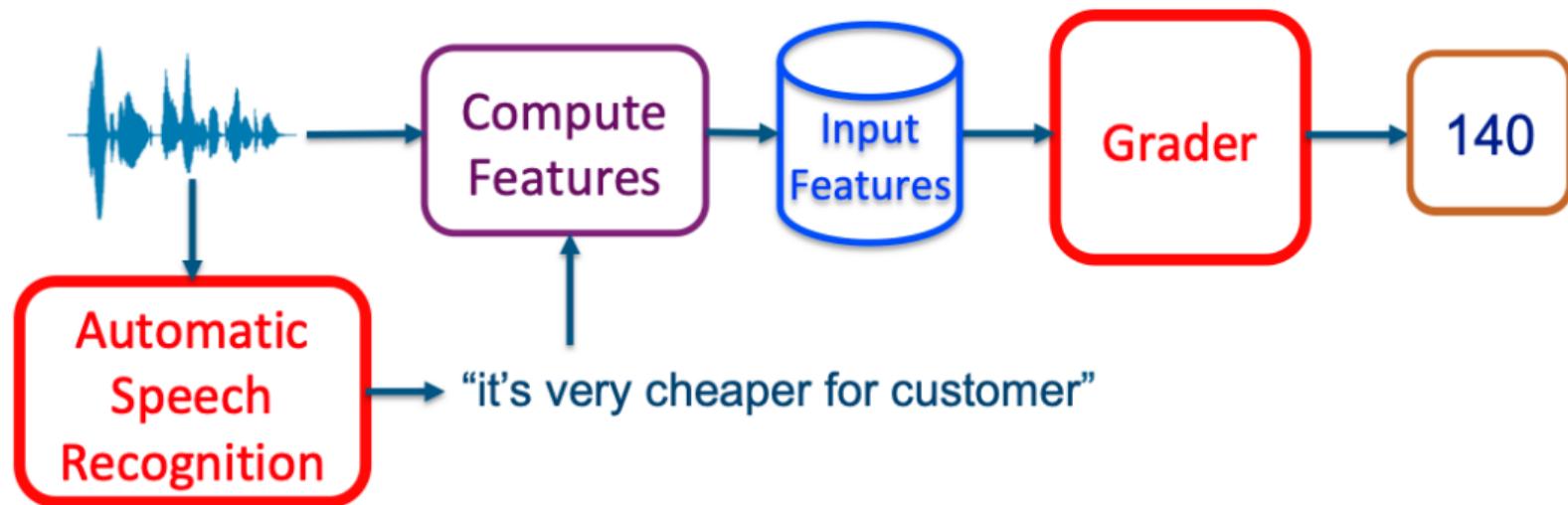
Mark Gales with Anton Ragni, Andrey Malinin  
and ALTA Speech Group

24th October 2019

# Computer says no: Irish vet fails oral English test needed to stay in Australia

**Louise Kennedy, a native English speaker with two degrees, says flawed technology is to blame**

# Spoken Language Assessment Pipeline

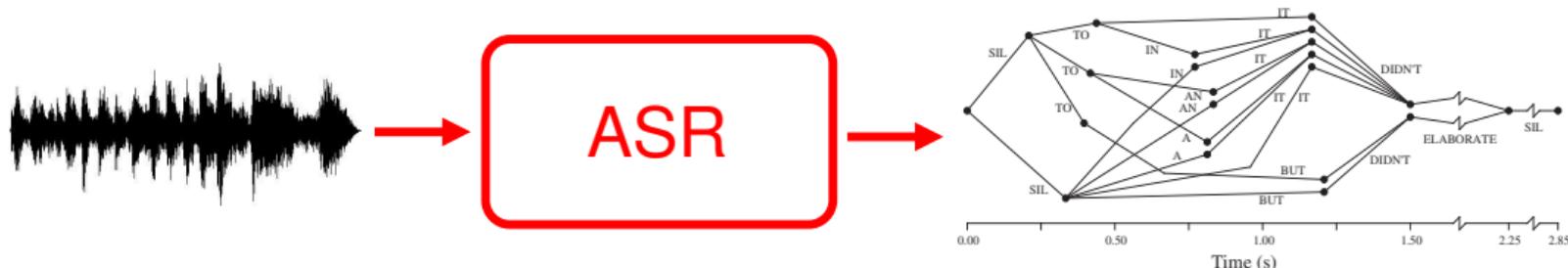


# ASR Confidence Scores

- Useful to know whether ASR output is correct
  - confidence scores supply this information
  - three forms of error: [substitutions](#), [deletions](#) and [insertions](#)

manual	AND	THESE	ARE		THE	FIMBLES
asr		THIS	ARE	TO	THE	FIMBLES
error	del	sub	—	ins	—	—
conf	—	0.4	0.8	0.3	0.9	0.9

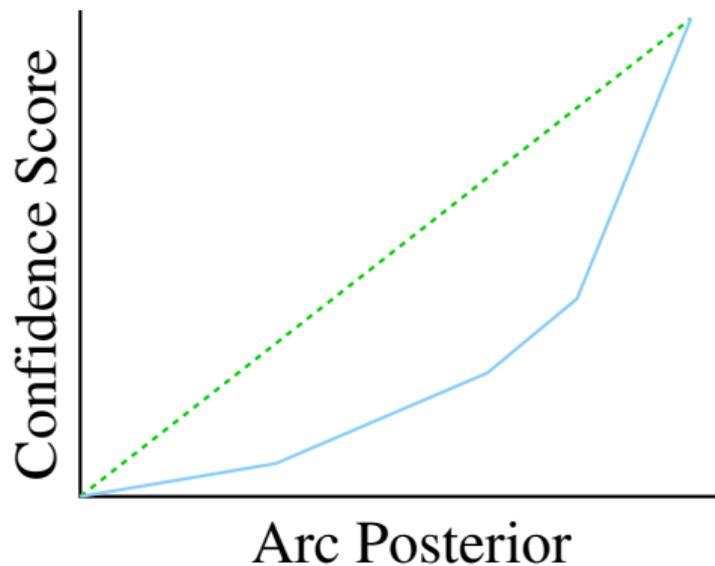
# Baseline Confidence Scores



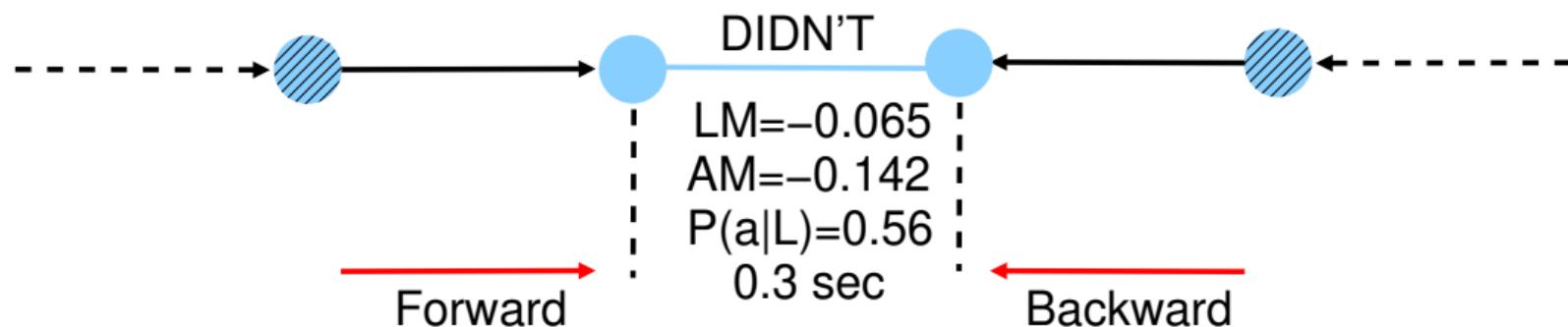
- Baseline confidence scores based on [arc posteriors](#)

$$p(\mathbf{q}_{1:T}, \mathbf{x}_{1:T}) = p_a(\mathbf{x}_{1:T} | \mathbf{q}_{1:T})^{\frac{1}{\gamma}} P_1(w_{1:L}); \quad P(a | \mathcal{L}) = \frac{\sum_{\mathbf{q}_{1:T} \in \mathcal{Q}_a} P(\mathbf{q}_{1:T}, \mathbf{x}_{1:T})}{p(\mathbf{x}_{1:T})}$$

- $\mathbf{q}_{1:T}$   $T$ -length state sequence for word sequence  $\omega$
- $\mathcal{Q}_a$  set of state sequences that pass through arc  $a$
- $\gamma$  is usually the LM scale factor
- does not alter 1-best (compared to scaling LM)



- Confidence scores often over-estimated
  - very simple normalisation approach



- Use more general sequence model
  - for 1-best  $w_{1:L} = w_1, \dots, w_L$
  - use information associated with each arc  $a_{1:L}$

$$P(w_i | a_{1:L}) = \mathcal{F}(a_i, \vec{a}_{1:i-1}, \overleftarrow{a}_{i+1:L})$$

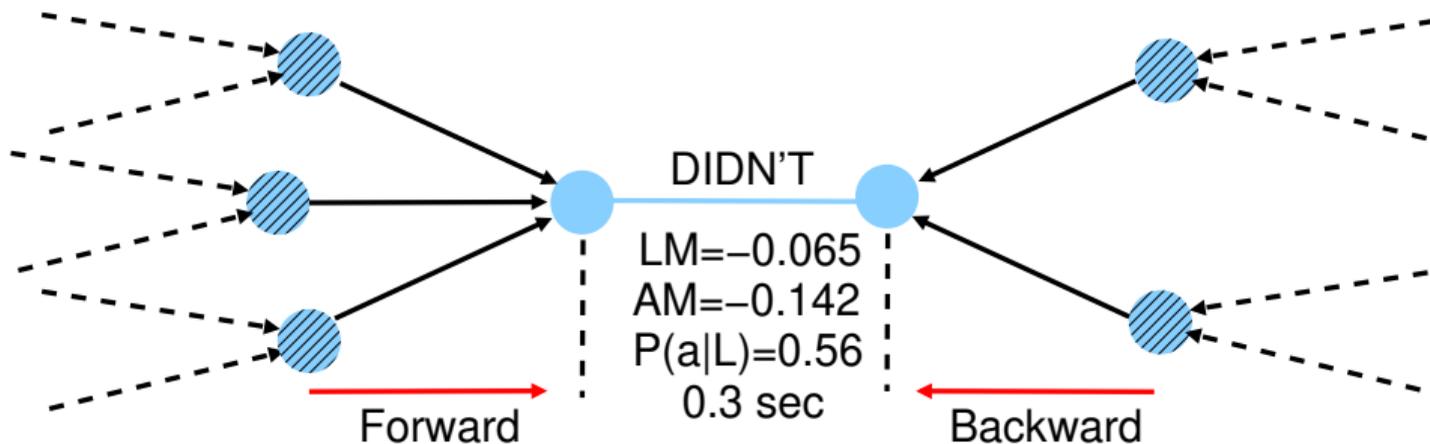
- Simple approach use recurrent neural networks

$$\vec{\mathbf{h}}_i = \mathcal{F}(\vec{\mathbf{h}}_{i-1}, \mathbf{a}_1); \quad \overleftarrow{\mathbf{h}}_i = \mathcal{F}(\overleftarrow{\mathbf{h}}_{i+1}, \mathbf{a}_1);$$
$$P(w_i | \mathbf{a}_{1:L}) = \mathcal{F}(\vec{\mathbf{h}}_i, \overleftarrow{\mathbf{h}}_i)$$

- Evaluation: Georgian (!) Conversation Telephone Speech
  - RNN-based on: posteriors, word ID and durations

System	NCE	AUC
Arc posteriors	-0.1978	0.9081
+ calibration	0.2755	0.9081
+ RNN	0.2911	0.9121

# Lattice Neural Network Based Confidence Scores [?, ?]

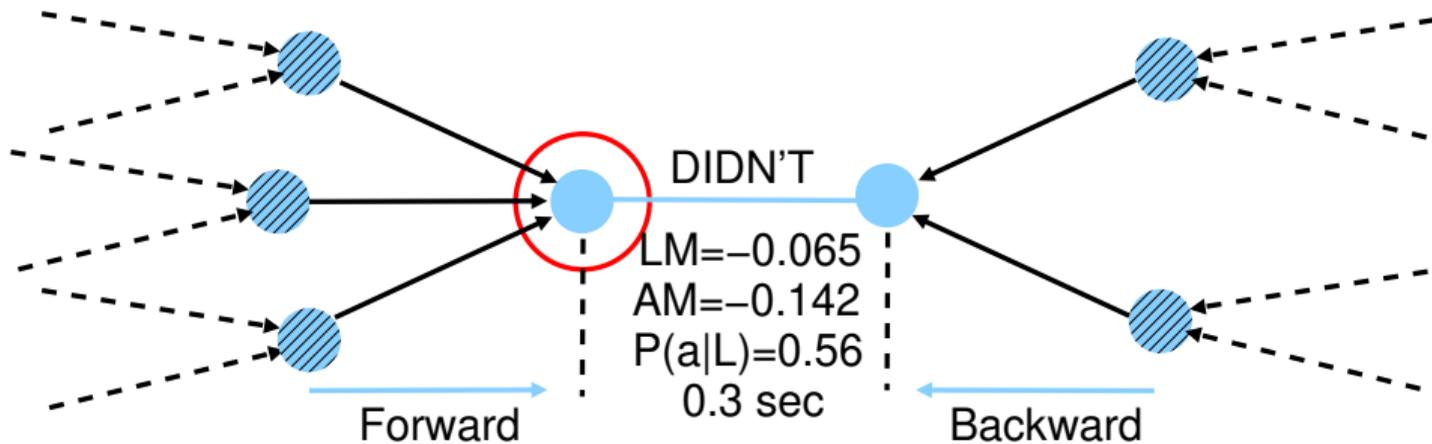


- Make use of complete lattice  $\mathcal{L}$

$$P(w_i|\mathcal{L}) = \mathcal{F}(a_i, \vec{Q}_{a_i}, \overleftarrow{Q}_{a_i})$$

- $\vec{Q}_{a_i}$  set of arcs in forward direction to  $a_i$
- $\overleftarrow{Q}_{a_i}$  set of arcs in backward direction to  $a_i$

# Lattice Neural Network Based Confidence Scores



- Use **attention** to merge arcs

$$\vec{h}_i = \text{attention} \left( \{ \vec{h}_j \}_{j \in \vec{N}_{a_i}}, a_i \right); \quad \overleftarrow{h}_i = \text{attention} \left( \{ \overleftarrow{h}_j \}_{j \in \overleftarrow{N}_{a_i}}, a_i \right);$$

$$P(w_i | a_{1:L}) = \mathcal{F}(\vec{h}_i, \overleftarrow{h}_i)$$

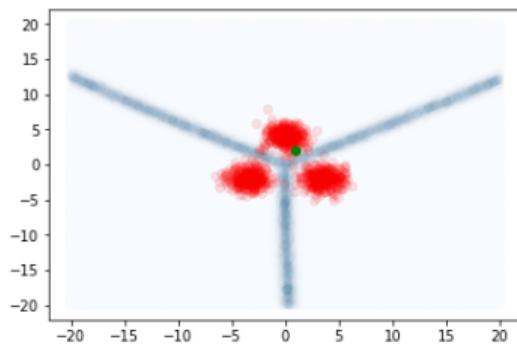


- Evaluation on a CTS task
  - RNN-based on: posteriors, word ID and durations
  - latticeRNN acts on [confusion networks](#)

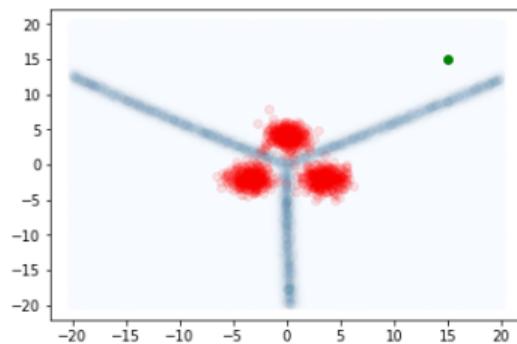
System	NCE	AUC
RNN	0.2911	0.9121
lattice-RNN	0.2934	0.9178
+ grapheme	0.3004	0.9231

# Prediction Uncertainty

# Sources of Uncertainty

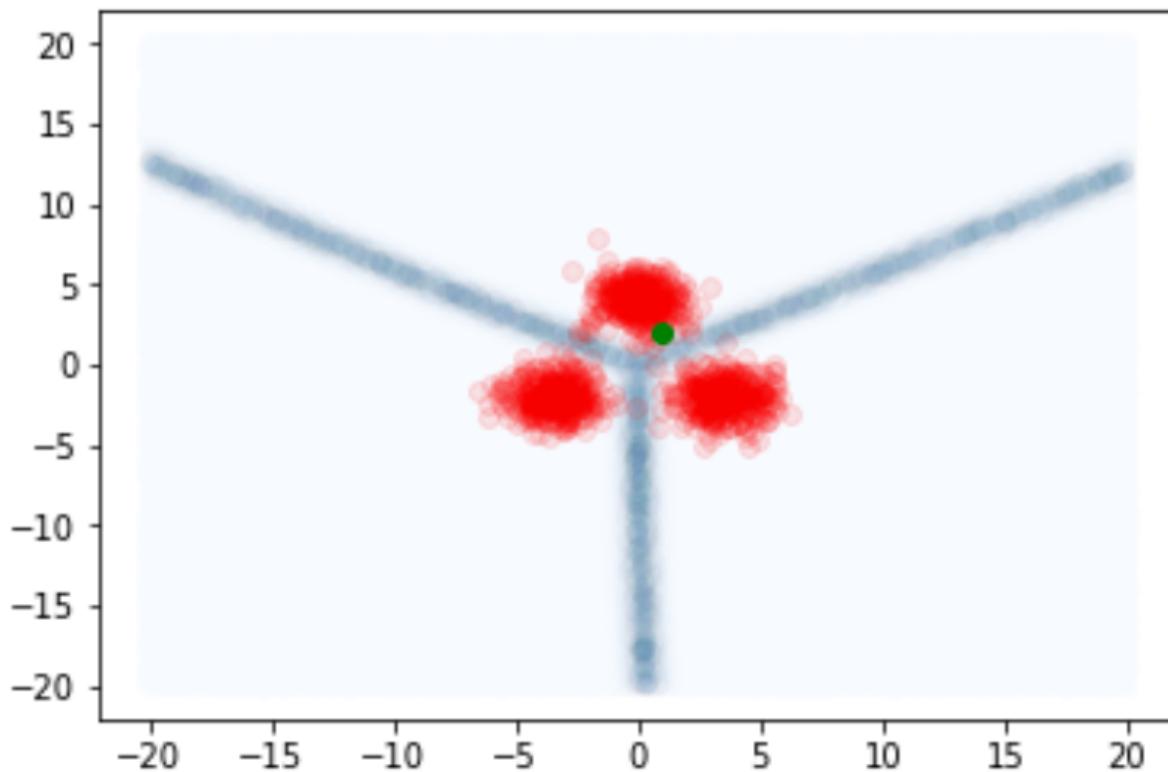


(a) Data Uncertainty



(b) Knowledge Uncertainty

# Data (Aleatoric) Uncertainty



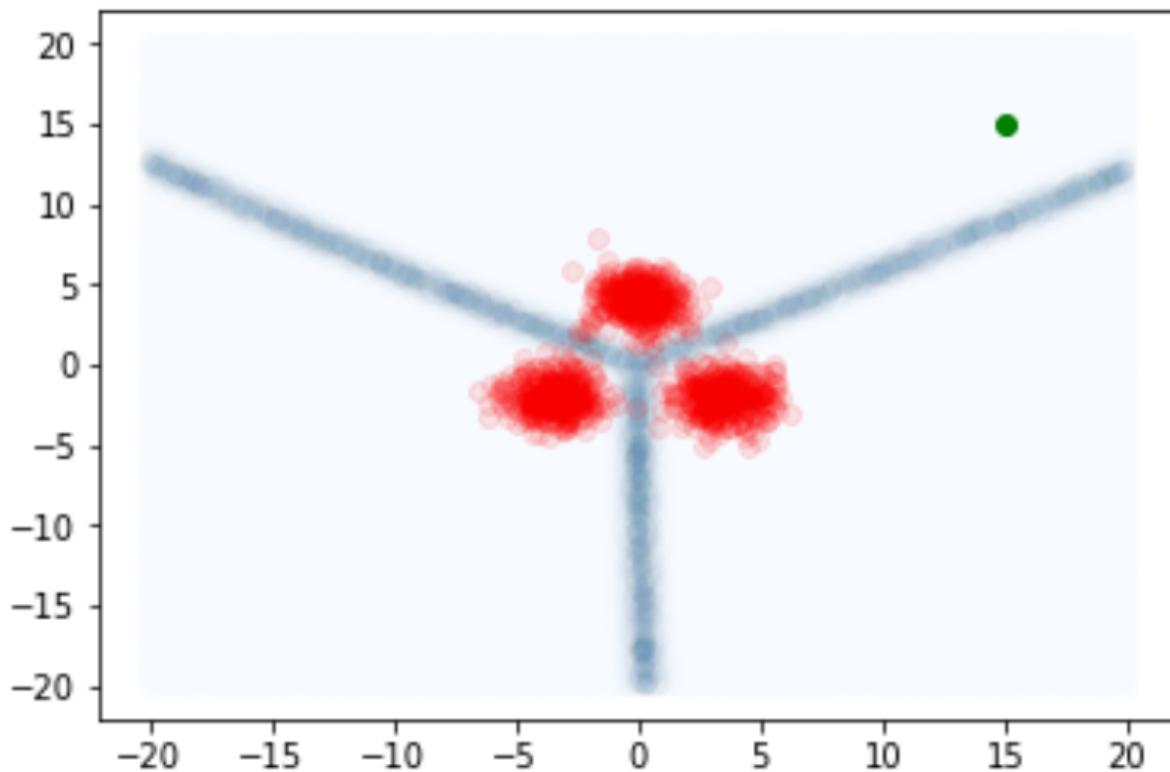
- Distinct Classes



- Overlapping Classes

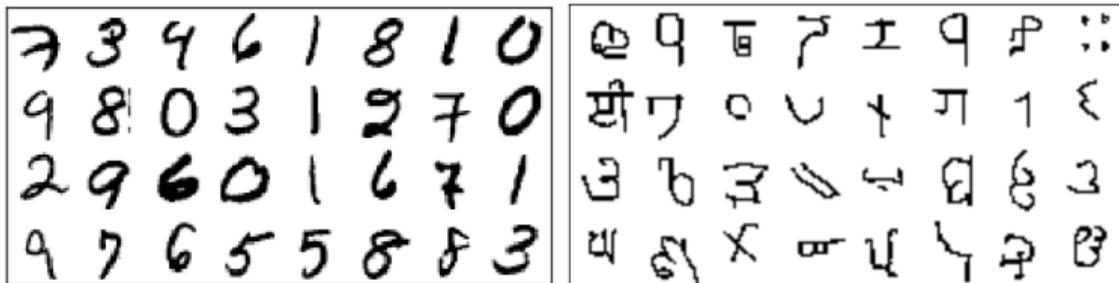


# Knowledge (Distributional/Epistemic) Uncertainty



# Knowledge Uncertainty

- Unseen classes



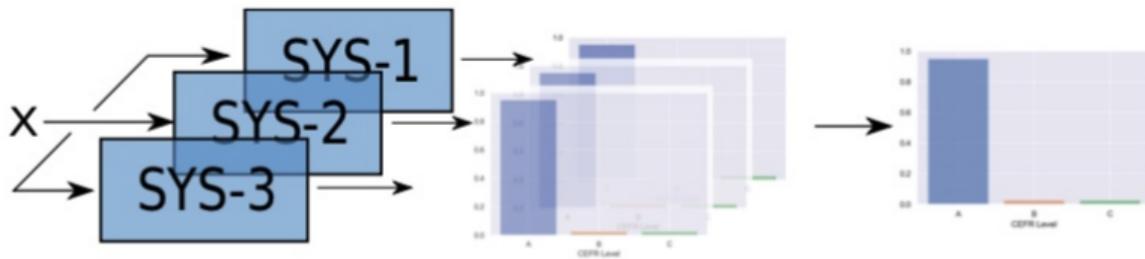
- Unseen variations of seen classes



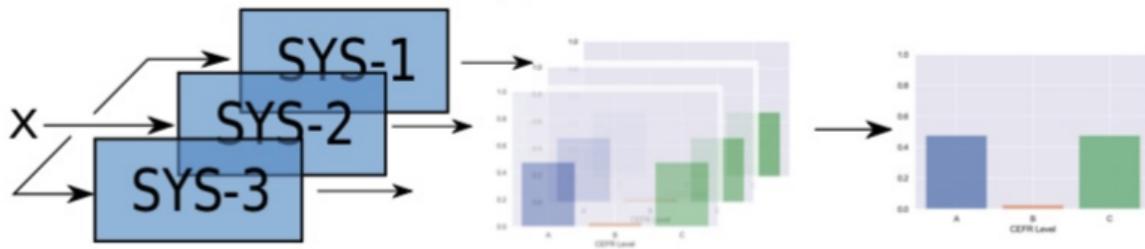
- Given training data  $\mathcal{D}$

$$\begin{aligned} P(y|\mathbf{x}^*, \boldsymbol{\theta}) &= \int P(y|\mathbf{x}^*, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} \\ &\approx \frac{1}{M} \sum_{i=1}^M P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(i)}); \quad \boldsymbol{\theta}^{(i)} \sim p(\boldsymbol{\theta}|\mathcal{D}) \end{aligned}$$

# Ensemble Approaches



(a) Confident



(b) Uncertain on decision boundary



- Simple reminder of Entropy

$$\mathcal{H}[\mathbb{P}(y|\mathbf{x}^*, \boldsymbol{\theta})] = - \sum_{c=1}^K \mathbb{P}(y = \omega_c | \mathbf{x}^*, \boldsymbol{\theta}) \log (\mathbb{P}(y = \omega_c | \mathbf{x}^*, \boldsymbol{\theta}))$$

- General attributes
  - high entropy: “flat” distribution, low confidence
  - low entropy: “peaky” distribution, high confidence
- Doesn't give information about source of uncertainty!

- Mutual Information

$$\underbrace{\mathcal{I}[y, \boldsymbol{\theta} | \mathbf{x}^*, \mathcal{D}]}_{\text{Knowledge Uncertainty}} = \underbrace{\mathcal{H}[\mathbb{E}_{\mathbf{p}(\boldsymbol{\theta} | \mathcal{D})}[\mathbf{p}(y | \mathbf{x}^*, \boldsymbol{\theta})]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{\mathbf{p}(\boldsymbol{\theta} | \mathcal{D})}[\mathcal{H}[\mathbf{p}(y | \mathbf{x}^*, \boldsymbol{\theta})]]}_{\text{Expected Data Uncertainty}}$$

- Total Variance

$$\underbrace{\mathbb{V}[y, \boldsymbol{\theta} | \mathbf{x}^*, \mathcal{D}]}_{\text{Total Variance}} = \underbrace{\mathbb{V}_{\mathbf{p}(\boldsymbol{\theta} | \mathcal{D})}[\mathbb{E}_{\mathbf{p}(y | \mathbf{x}^*, \boldsymbol{\theta})}[y]]}_{\text{Mean Variance}} + \underbrace{\mathbb{E}_{\mathbf{p}(\boldsymbol{\theta} | \mathcal{D})}[\mathbb{V}_{\mathbf{p}(y | \mathbf{x}^*, \boldsymbol{\theta})}[y]]}_{\text{Expected Data Variance}}$$

- Expected (Pairwise) KL-Divergence

$$\text{KL}[y, \boldsymbol{\theta} | \mathbf{x}^*, \mathcal{D}] = \mathbb{E}_{\mathbf{p}(\boldsymbol{\theta} | \mathcal{D}), \mathbf{p}(\tilde{\boldsymbol{\theta}} | \mathcal{D})} [\text{KL}[\mathbf{p}(y | \mathbf{x}^*, \boldsymbol{\theta}) || \mathbf{p}(y | \mathbf{x}^*, \tilde{\boldsymbol{\theta}})]]$$

- Deep learning approaches often use 10,000,000+ parameters

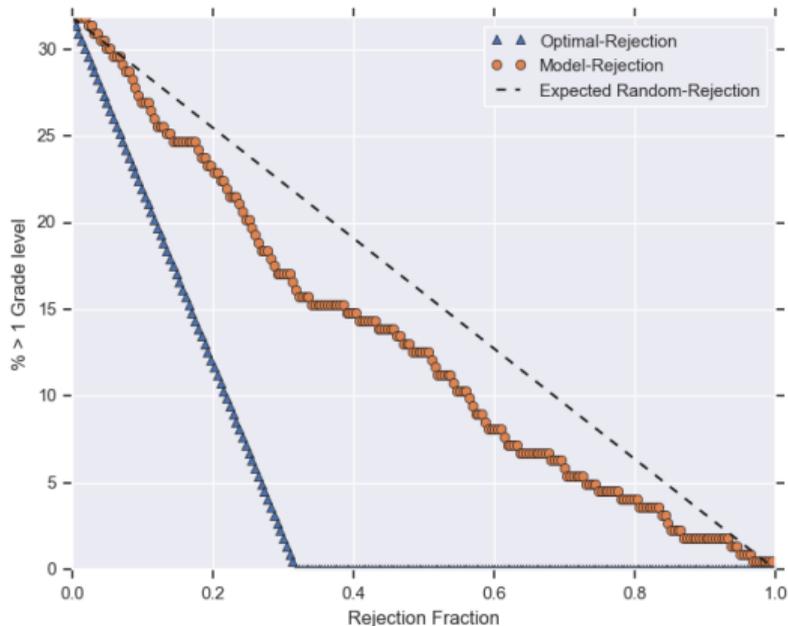
Modelling  $p(\theta|\mathcal{D})$  challenging

- use variational approximations
- Monte-Carlo methods
- non-Bayesian approaches e.g. random network initialisation

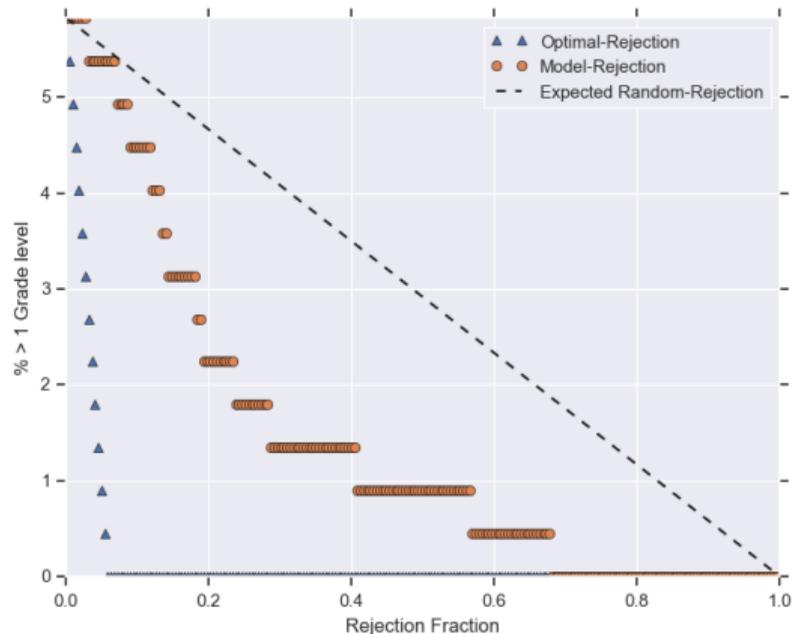
# Computer says no: Irish vet fails oral English test needed to stay in Australia

**Louise Kennedy, a native English speaker with two degrees, says flawed technology is to blame**

# Grader Uncertainty: Ensemble-Based



Within 0.5 CEFR-Level



Within 1.0 CEFR-Level

# Prior Networks

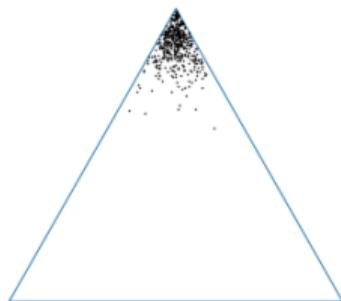
- Ensembles compute/memory intensive (scales linearly)
  - challenging to [guarantee](#) performance for outliers

- Ensembles compute/memory intensive (scales linearly)
  - challenging to **guarantee** performance for outliers
- Possible to compress ensemble to a single model:
  - **Ensemble Distillation**: standard compression approach

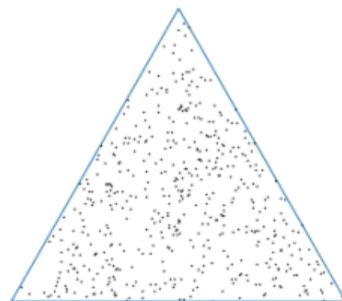
$$\hat{\theta} = \arg \max_{\theta} \left\{ \text{KL} \left( \frac{1}{M} \sum_{i=1}^M P(y|\mathbf{x}^*, \theta^{(i)}) \parallel P(y|\mathbf{x}^*, \theta) \right) \right\}$$

- models average distribution - loses diversity of ensemble
- **Ensemble Distribution Distillation**: model ensemble diversity
  - maintains diversity of the ensemble

- Ensemble  $\{P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(i)})\}_{i=1}^M$  can be visualised on a simplex



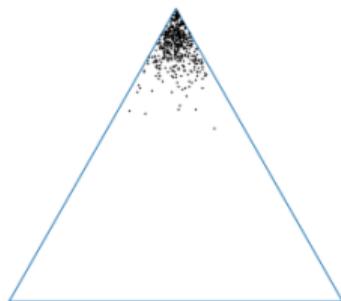
(a) In domain  $\mathbf{x}^*$



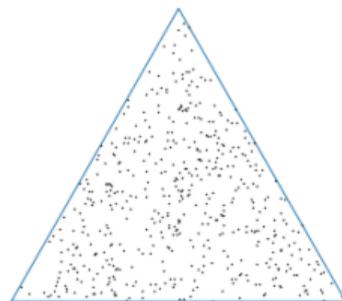
(b) Out-of-domain  $\mathbf{x}^*$

- ensemble samples from a **distribution over distributions**

- Ensemble  $\{P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(i)})\}_{i=1}^M$  can be visualised on a simplex



(a) In domain  $\mathbf{x}^*$



(b) Out-of-domain  $\mathbf{x}^*$

- ensemble samples from a **distribution over distributions**
- Only** need to model desired distribution
  - should allow **explicit** control over diversity

- A Prior Network predicts parameters of **Dirichlet Distribution**

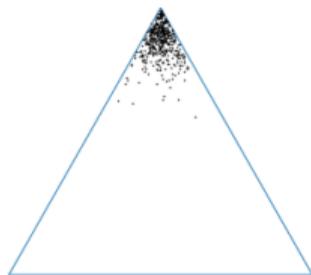
$$p(\boldsymbol{\mu}|\mathbf{x}^*; \hat{\boldsymbol{\theta}}) = \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}), \quad \boldsymbol{\alpha} = \mathbf{f}(\mathbf{x}^*; \hat{\boldsymbol{\theta}})$$

where

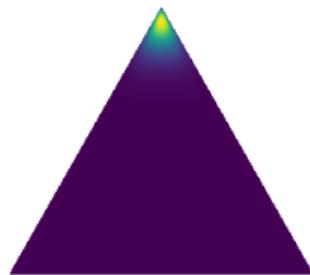
$$\boldsymbol{\mu} = \begin{bmatrix} P(y = \omega_1 | \mathbf{x}^*) \\ P(y = \omega_2 | \mathbf{x}^*) \\ \vdots \\ P(y = \omega_K | \mathbf{x}^*) \end{bmatrix}$$

- Dirichlet Distribution  $\rightarrow$  Distribution over simplex
  - Conjugate prior to categorical distribution
  - Convenient properties  $\rightarrow$  analytically tractable

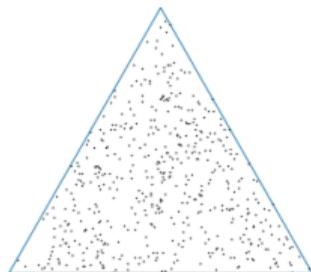
# Distribution over Distributions



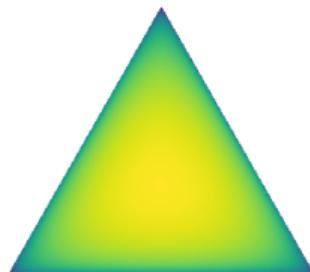
(a)  $\{\mu^{(i)}\}_{i=1}^M$



(b)  $p(\mu|x^*, \mathcal{D})$

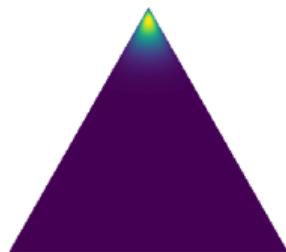


(c)  $\{\mu^{(i)}\}_{i=1}^M$

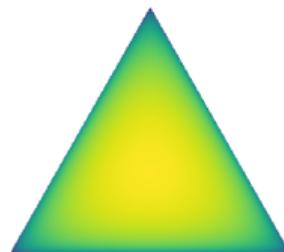


(d)  $p(\mu|x^*, \mathcal{D})$

$$\mathcal{L}(\theta, \mathcal{D}) = \underbrace{\mathcal{L}_{in}(\theta, \mathcal{D}_{trn})}_{\text{In Domain Loss}} + \gamma \cdot \underbrace{\mathcal{L}_{out}(\theta, \mathcal{D}_{out})}_{\text{OOD Loss}}$$



(a) In-Domain Target



(b) OOD Target

- Explicitly train the form of the Dirichlet distributions
  - **but** requires selection/generation of out-of-distribution data

# Target Dirichlet Parameters [?]

- Train network to predict **appropriate** distribution:
  - map  $y^{(i)} \rightarrow \beta^{(i)}$ : should yield **correct class** - minimise

$$\mathcal{L}(\theta; \mathcal{D}) = \sum_{i=1}^N \text{KL} \left( p(\mu | \beta^{(i)}) \parallel p(\mu | \mathbf{x}^{(i)}; \theta) \right)$$

- Consider setting  $\beta^{(i)}$  as follows  $\rightarrow$

$$\beta_k^{(i)} = \begin{cases} \beta + 1 & \text{if } y^{(i)} = \omega_k \\ 1 & \text{if } y^{(i)} \neq \omega_k \end{cases}$$

- if  $\beta$  is large  $\rightarrow$ : high confidence
  - if  $\beta$  is low  $\rightarrow$ : low confidence
  - If  $\beta$  is zero  $\rightarrow$ : flat (uniform) distribution
- Reverse-KL yields better results (see paper for reasons)

# Out-of-Distribution Detection: Image Tasks

- Use CIFAR-100 for out-of-distribution (OOD) training data
  - evaluate performance in detecting OOD test samples
  - metric AUC (average 10 randomly initialised models  $\pm 2\sigma$ )

Model	CIFAR-10		
	SVHN	LSUN	TinyImageNet
Ensemble	89.5 $\pm$ NA	93.2 $\pm$ NA	90.3 $\pm$ NA
Prior Network	<b>98.2</b> $\pm 1.1$	<b>95.7</b> $\pm 0.9$	<b>95.7</b> $\pm 0.7$

# Conclusions

# Uncertainty: Knowing What You Don't Know

- Uncertainty important for deploying machine learning
  - systems tend to be overly confident

# Uncertainty: Knowing What You Don't Know

- Uncertainty important for deploying machine learning
  - systems tend to be overly confident
- Knowing the cause of uncertainty useful
  - allows different actions to be taken to address uncertainty
  - applications: active learning, uncertainty for RL, ...

# Uncertainty: Knowing What You Don't Know

- Uncertainty important for deploying machine learning
  - systems tend to be overly confident
- Knowing the cause of uncertainty useful
  - allows different actions to be taken to address uncertainty
  - applications: active learning, uncertainty for RL, ...
- It's hard!
  - humans aren't too good at it either