



A Confidence-based Acquisition Model for Self-supervised Active Learning and Label Correction

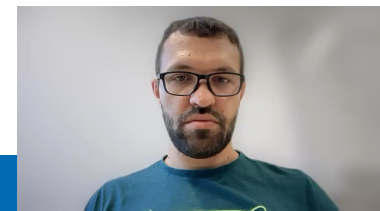
Carel van Niekerk, Christian Geishauser, Michael Heck, Shutong Feng, Hsien-chin Lin, Nurul Lubis, Benjamin Ruppik, Renato Vukovic, and Milica Gašić



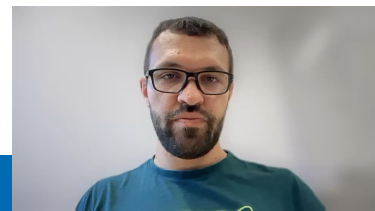
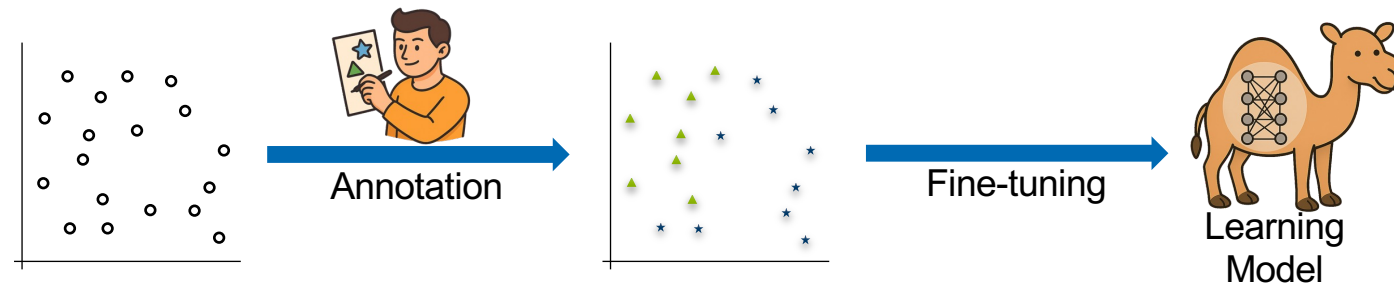
Contact Me



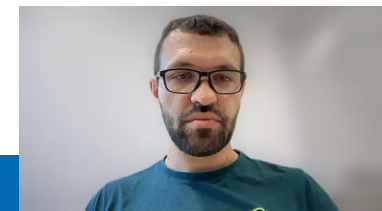
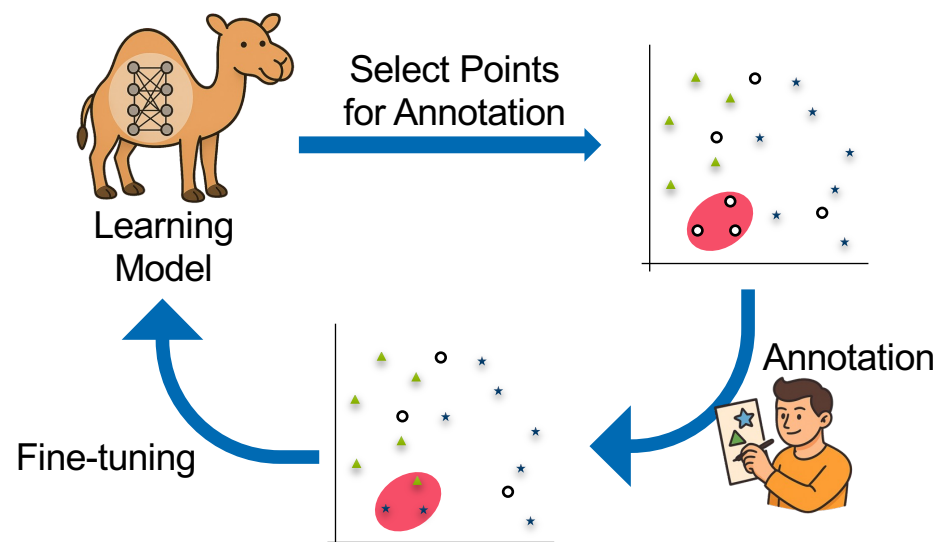
Implementation



- In **Passive Learning** **large annotated corpora** are collected for tasks such as machine translation, dialogue modelling, etc.

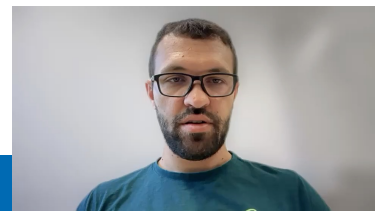
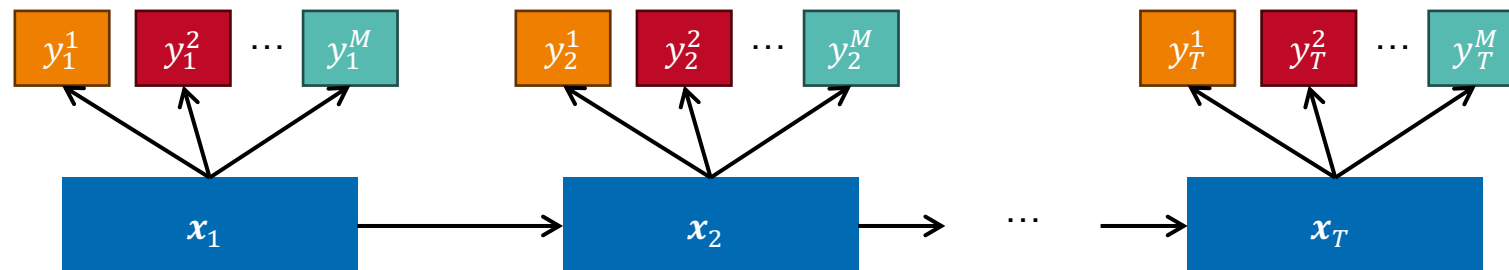


- In **Passive Learning** **large annotated corpora** are collected for tasks such as machine translation, dialogue modelling, etc.
- In **Active Learning**, the learning **model selects the most beneficial datapoints** to learn, reducing the annotation effort.



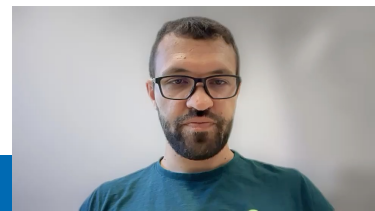
Sequential Multi-Output Problem

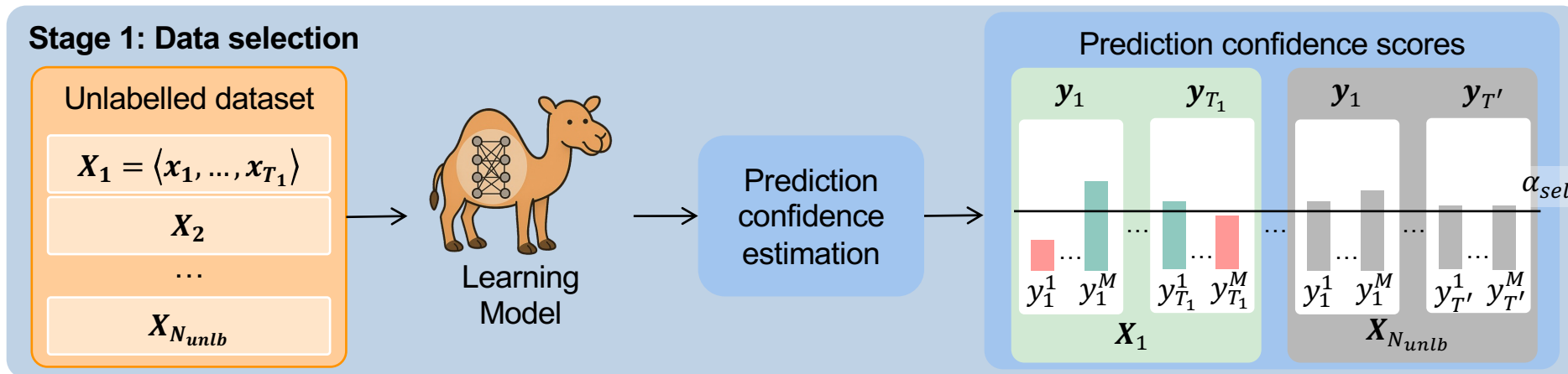
- **Sequential Multi-Output** problems require a label at each timestep for each output category.
- Expert labels can be very **expensive** and crowd labels very **noisy**.

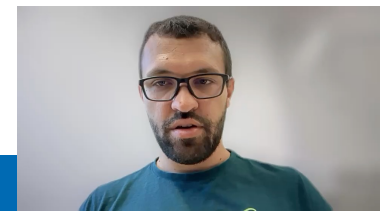
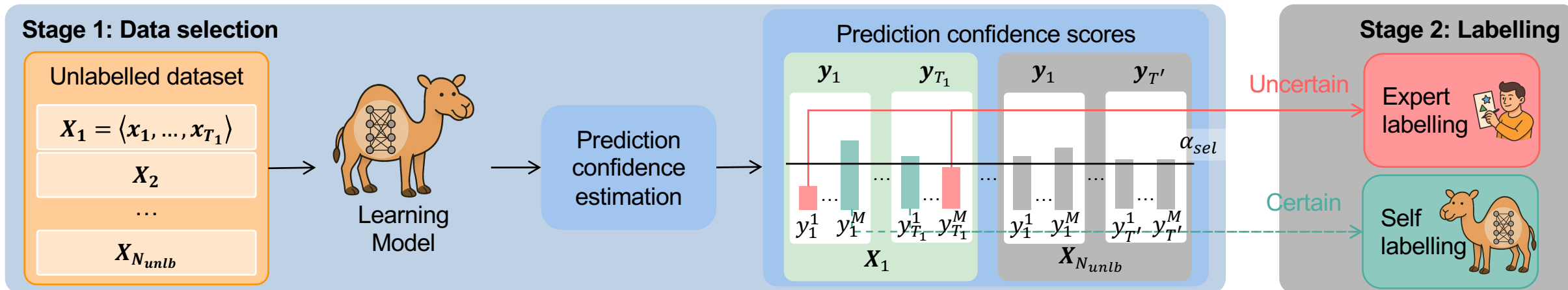


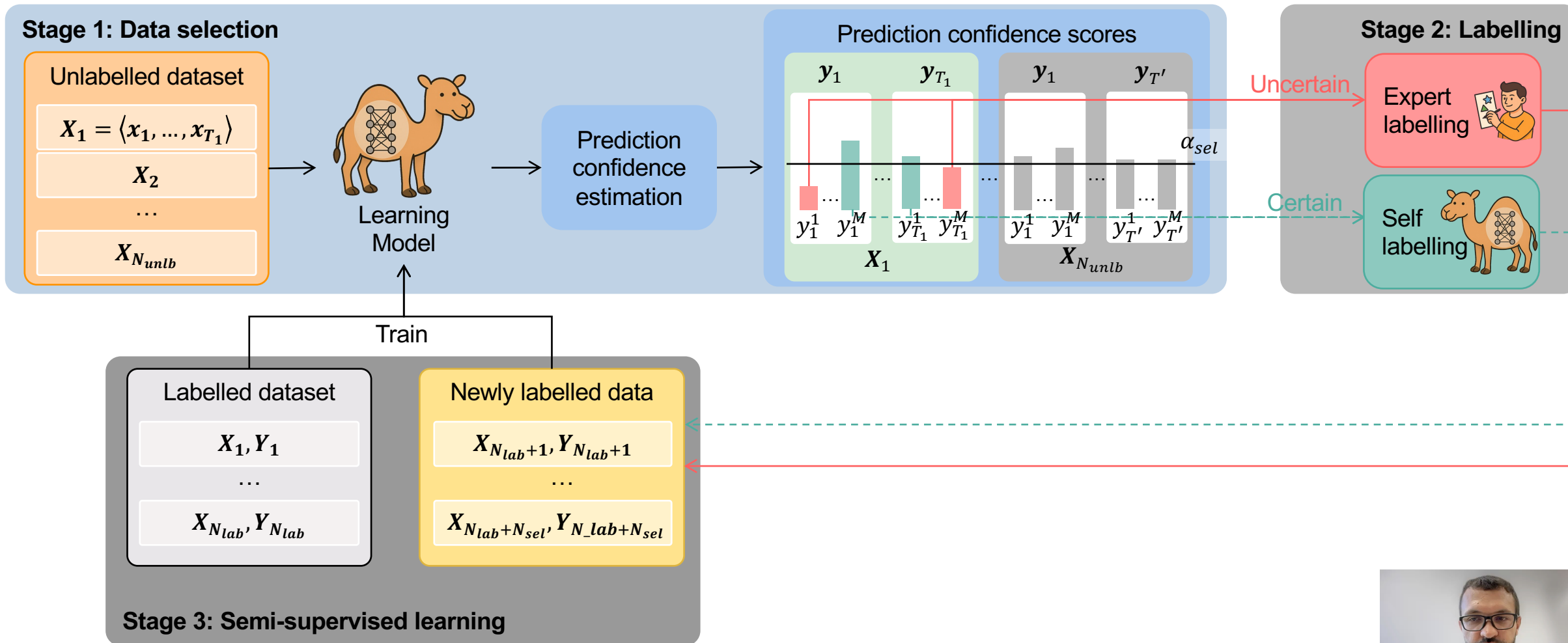
■ CAMEL:

- Confidence-based Acquisition Model
- for Efficient self-supervised active Learning

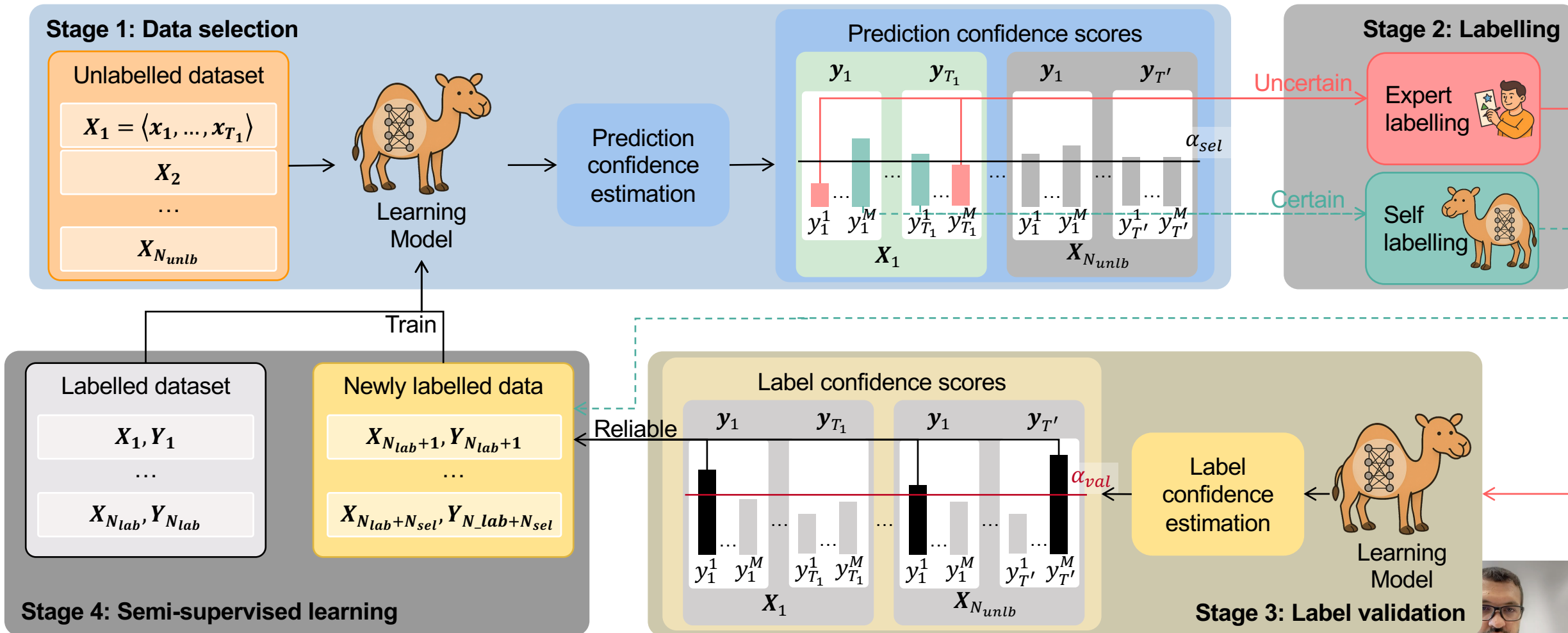






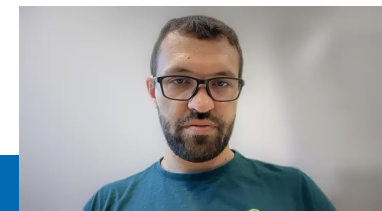
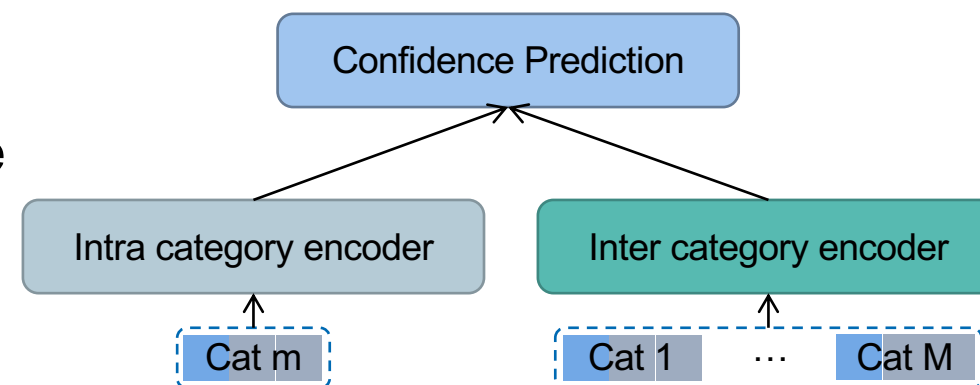


CAMELL – CAMEL with Label validation



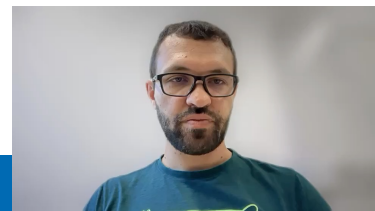
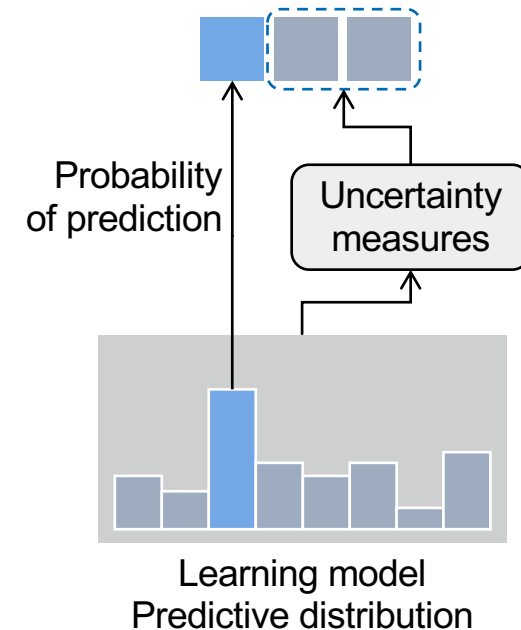
The Model

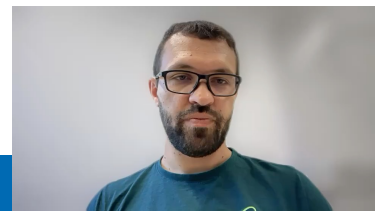
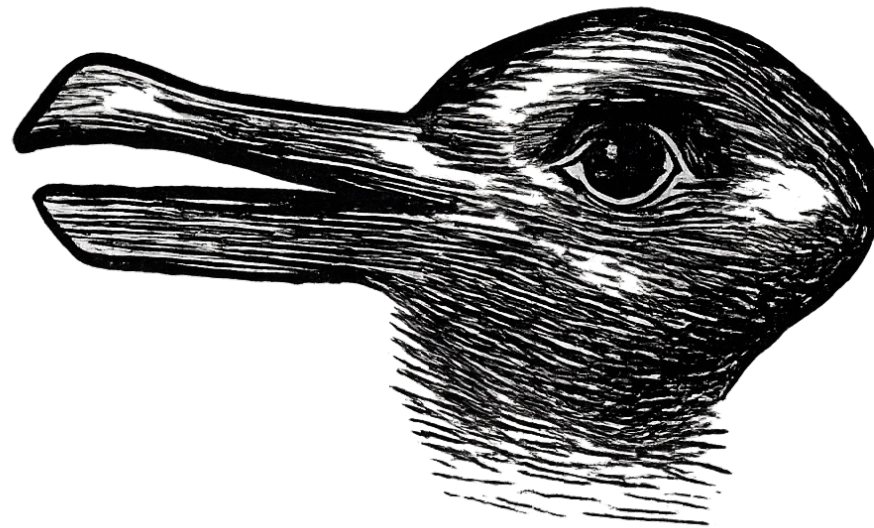
- Incorporates intra-category features to capture category specific uncertainty.
- Incorporates inter-category features to capture the correlation between categories.
- The combined intra- and inter-category encodings used for predicting the confidence.
- Objective: Predict whether the prediction of the learning model is correct.

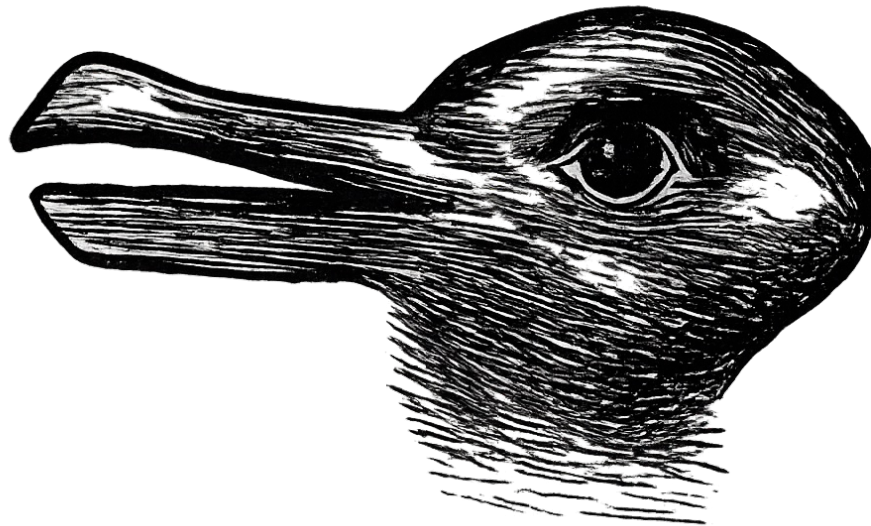


The uncertainty features

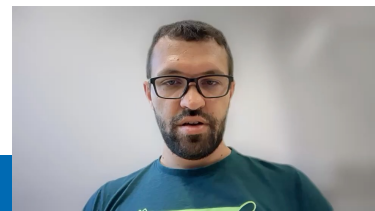
- Probability of the prediction / label.
- Uncertainty features extracted from the predictive distribution of the learning and noisy models:
 - Total Uncertainty (Entropy)
 - Knowledge Uncertainty

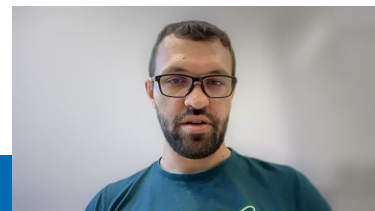
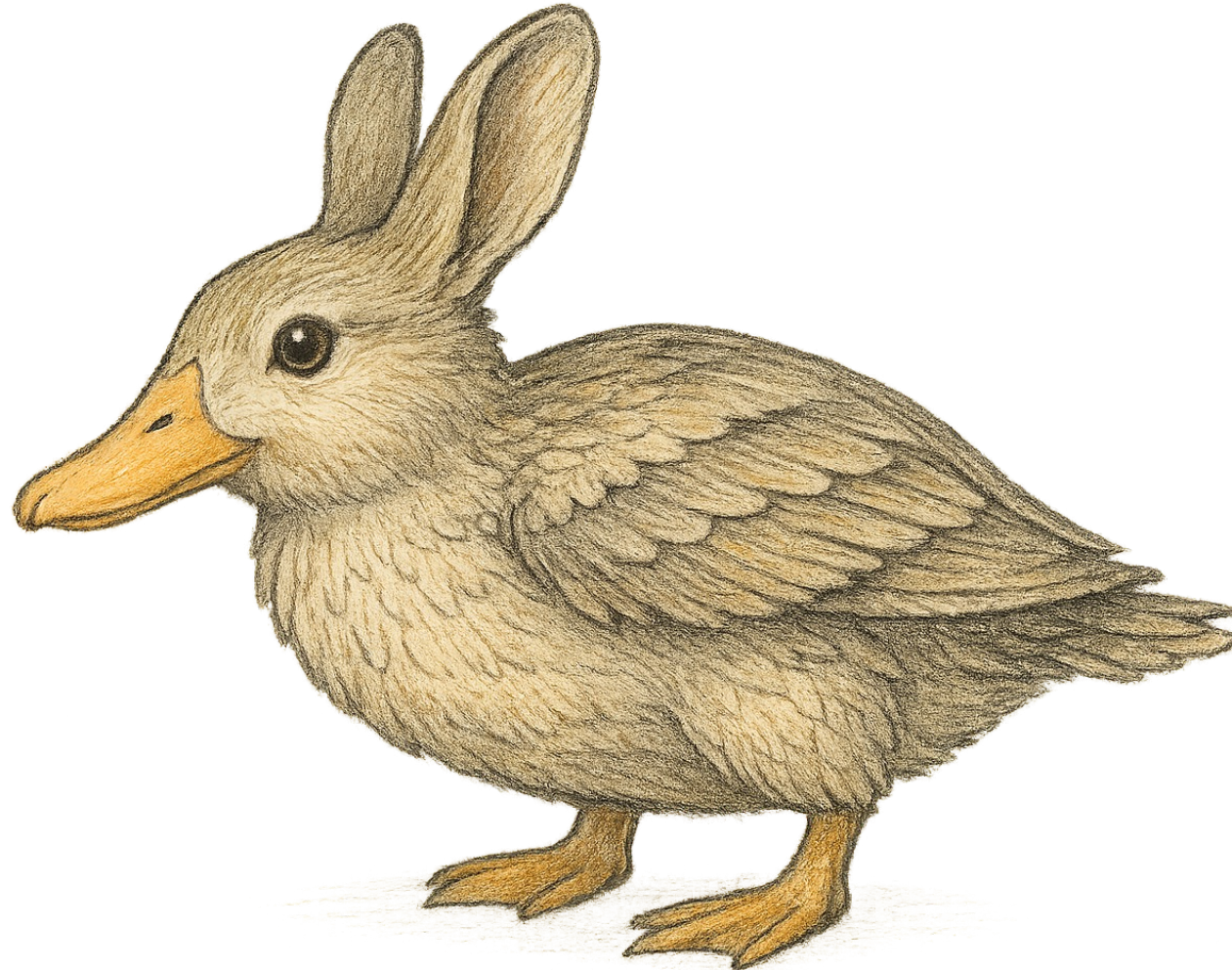


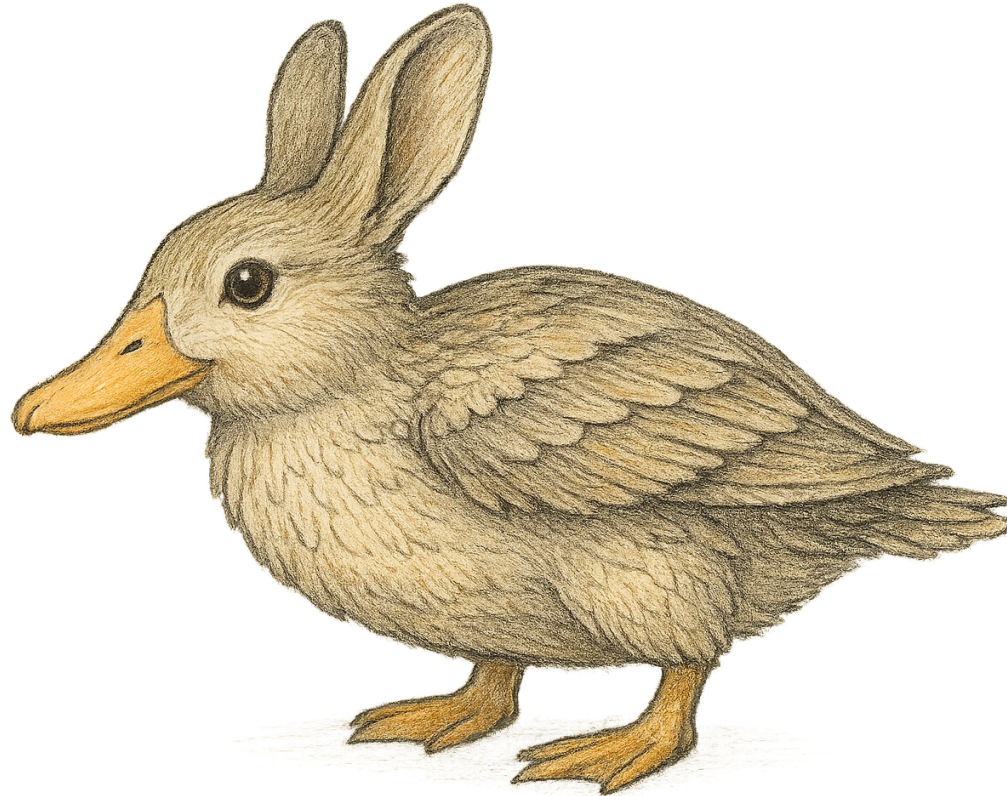




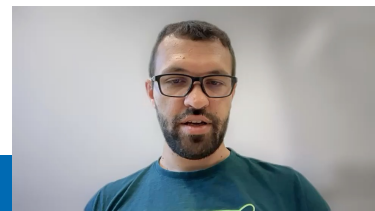
Duck	Rabbit
0.5	0.5

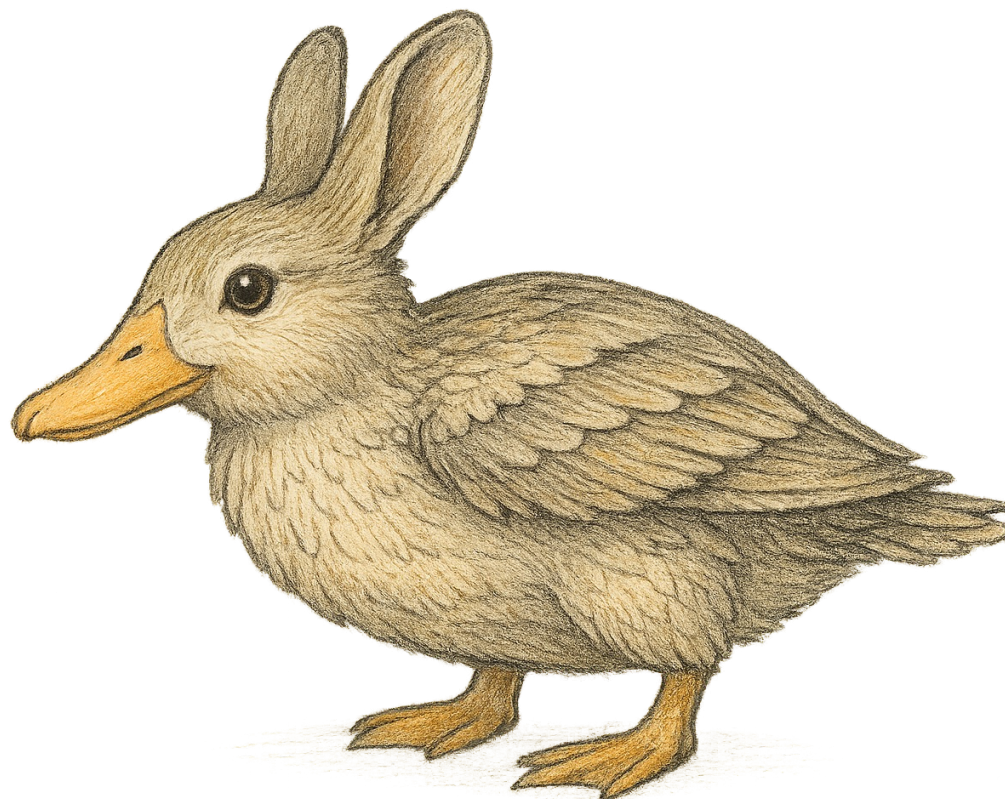




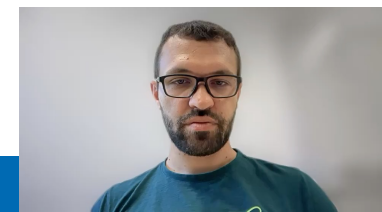


Duck	Rabbit
0.7	0.3

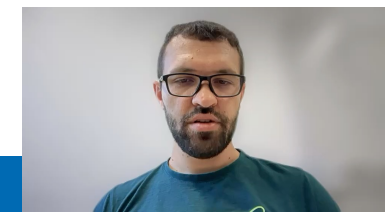
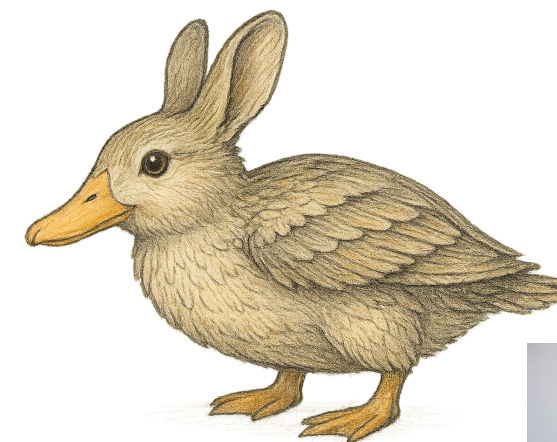
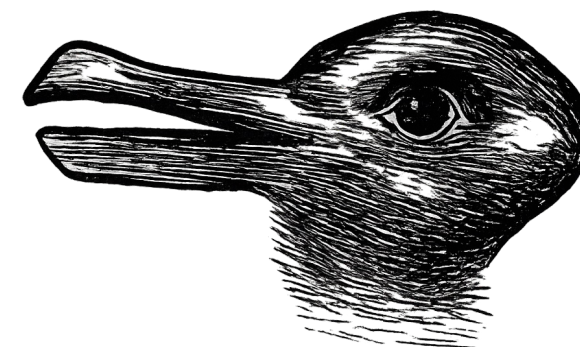




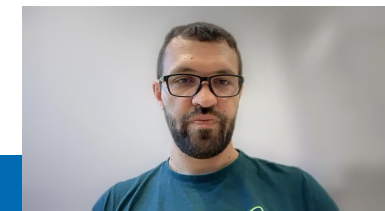
Duck	Rabbit	Dubbit
0.08	0.02	0.9



- Uncertainty in Machine Learning models stem from two main sources:
- **Data** Uncertainty: Uncertainty which stems from **ambiguity in the data**.
- **Knowledge** Uncertainty: Uncertainty which stems from a **lack of knowledge** within the model.

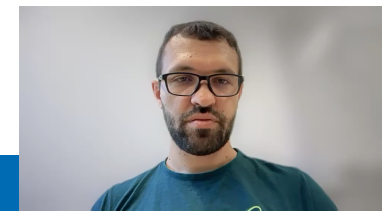
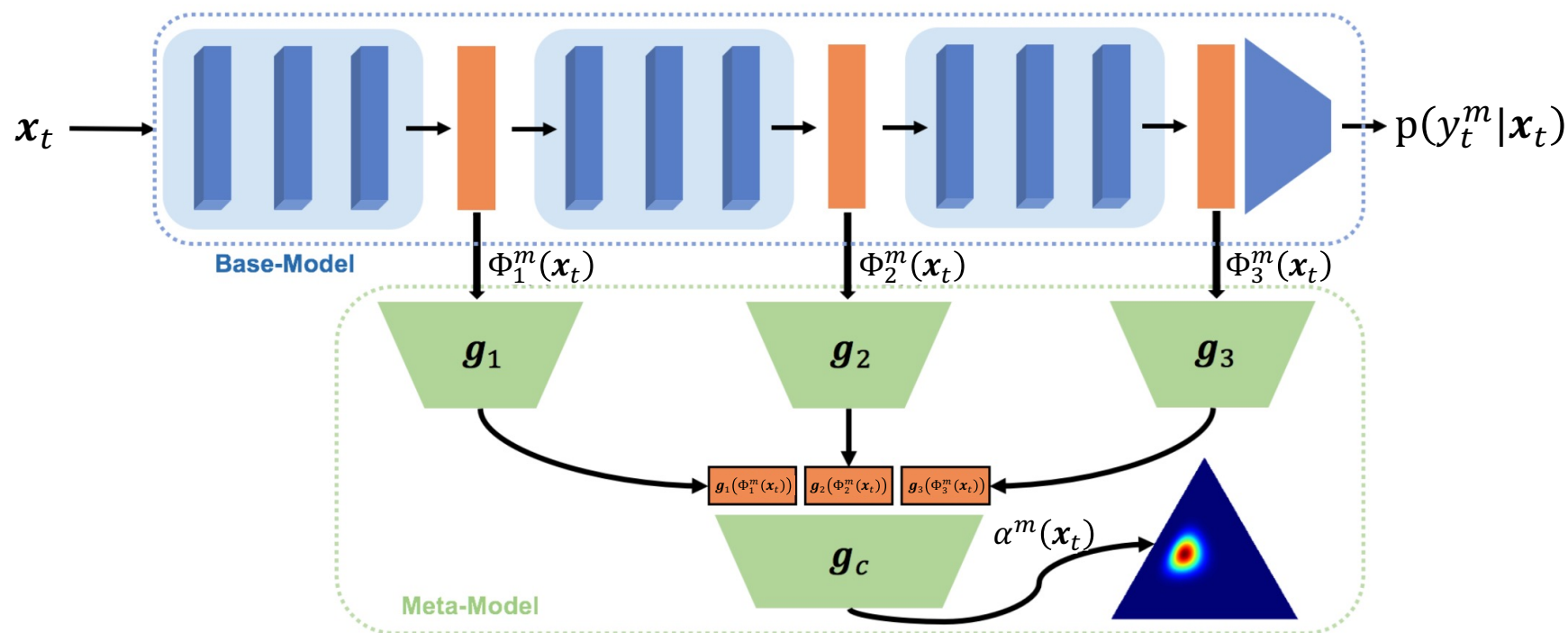


- The **source of uncertainty** can be distinguished using an ensemble of models.
 - The **entropy** in the predictive distribution is the **total uncertainty**.
 - The **variation** between ensemble member predictions is the **knowledge uncertainty**.
 - Data uncertainty is the difference between these two.
 - This is **computationally expensive** especially for language models.
- Alternative: Learning a **Dirichlet distribution**, $\text{Dirichlet}(\alpha)$, over the probability simplex.
 - The **mean predictive distribution** provides a measure of **total uncertainty**.
 - The **variation** under this distribution provides a measure of **knowledge uncertainty**.



Uncertainty Estimation

- The **variation between embeddings** at different layers allows accurate post-hoc uncertainty learning.

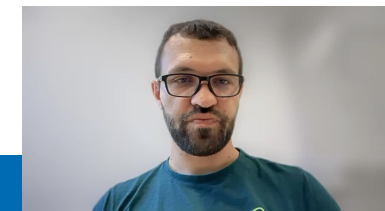


Uncertainty Estimation Objective

- ELBO Objective:

$$\mathcal{L}(\boldsymbol{\theta}^{(\text{meta})}; \mathcal{D})$$

- Expected **Likelihood**: $= \mathbb{E}_{p(\mathbf{x}, y | \mathcal{D})} \left[\mathbb{E}_{p(\boldsymbol{\pi} | \mathbf{x}, \boldsymbol{\theta}^{(\text{meta})})} [-\log p(y | \boldsymbol{\pi})] \right]$
- **KL Penalty**: $+ \lambda \mathbb{E}_{p(\mathbf{x}, y | \mathcal{D})} \left[D_{KL} [p(\boldsymbol{\pi} | \mathbf{x}, \boldsymbol{\theta}^{(\text{meta})}) || p(\boldsymbol{\pi} | \boldsymbol{\beta})] \right]$
- Challenge: We require an **informative** prior $\text{Dirichlet}(\boldsymbol{\beta})$.



- Train an ensemble of E models on a **small subset** of the training data.
- Using predictions from this ensemble to produce the prior using **Sterlings Approximation**:

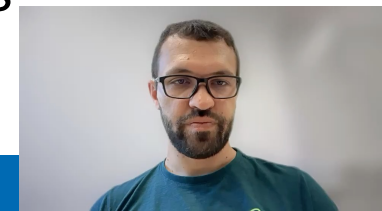
$$\beta = \beta_0(\mathbf{x})\tilde{\pi}(\mathbf{x}), \text{ where}$$

$$\tilde{\pi}(\mathbf{x}) = \frac{1}{E} \sum_{e=1}^E \pi_k^{(e)}(\mathbf{x}) \quad (\text{Mean – Total Uncertainty})$$

$$\beta_0(\mathbf{x}) = \frac{K-1}{2 \sum_{k=1}^K \tilde{\pi}_k(\mathbf{x}) d_k(\mathbf{x})}, \text{ where}$$

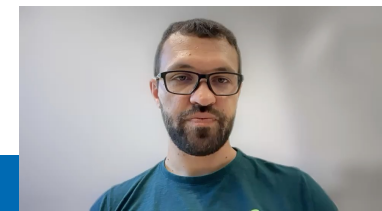
$$d_k(\mathbf{x}) = \log \tilde{\pi}_k(\mathbf{x}) - \frac{1}{E} \sum_{e=1}^E \log \pi_k^{(e)}(\mathbf{x}) \quad (\text{Variation – Knowledge Uncertainty})$$

- This prior is used in the first active learning step.
- After this the predicted posterior of the previous active learning step is used as a prior in order to **update the beliefs of the model**.

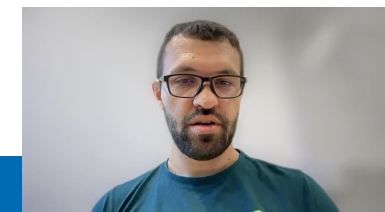
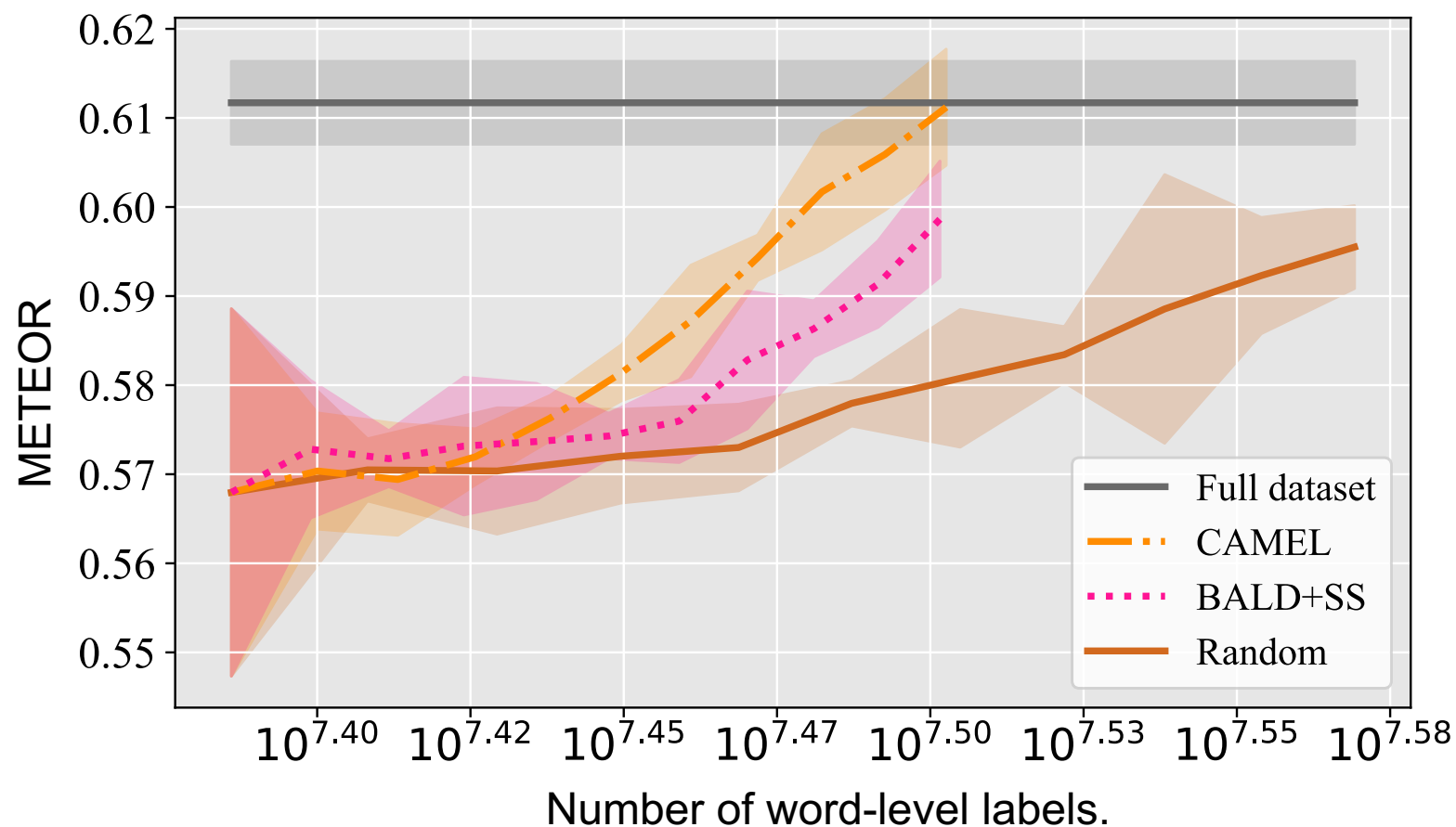


Machine Translation

- Model: Ensemble T5-small encoder-decoder transformer.
- Uncertainty Estimation:
 - Total Uncertainty: Entropy of the predictive distribution.
 - Knowledge Uncertainty: Mutual information between predictive distribution and ensemble members.
- Dataset: WMT17 DE-EN for German to English translation.
- Confidence Estimation Model - Simplified model for the single category sequential task.



Machine Translation



Dialogue Belief Tracking

■ Model

- Ensemble-SetSUMBT
- Meta-Uncertainty SetSUMBT

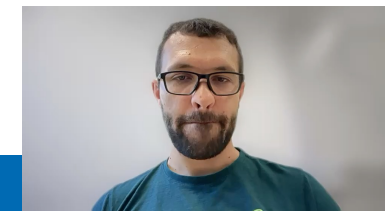
■ Uncertainty Estimation:

- Total Uncertainty: Entropy of the predictive distribution or Entropy within the Dirichlet distribution.
- Knowledge Uncertainty: Mutual information between predictive distribution and ensemble members or knowledge uncertainty estimate using the Dirichlet distribution.

■ Datasets:

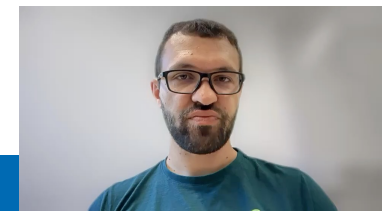
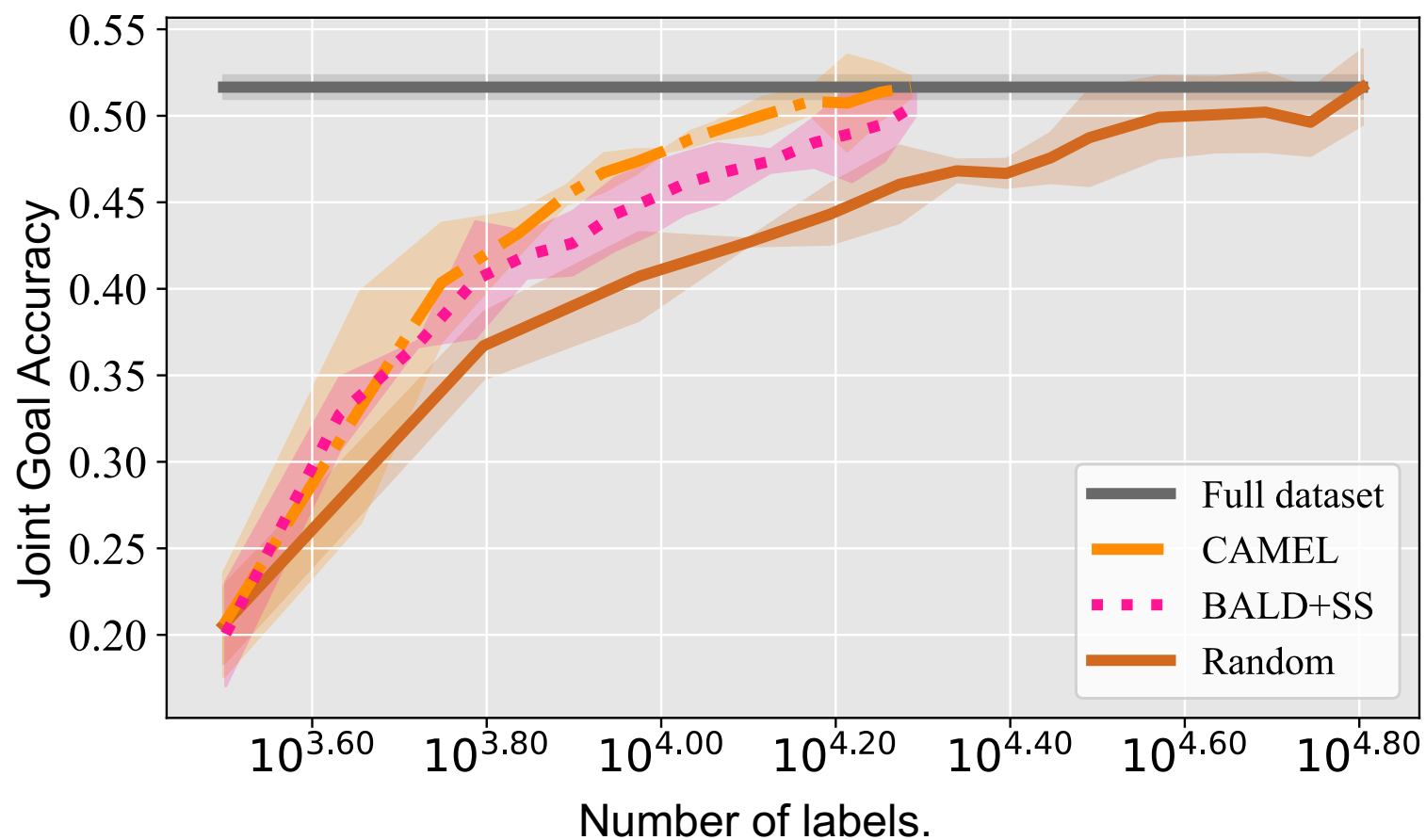
- MultiWOZ 2.1 (Noisy version)
- MultiWOZ 2.4 (Cleaned version)

■ Metric: Joint goal accuracy.

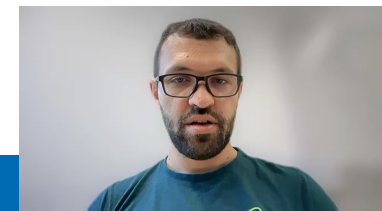
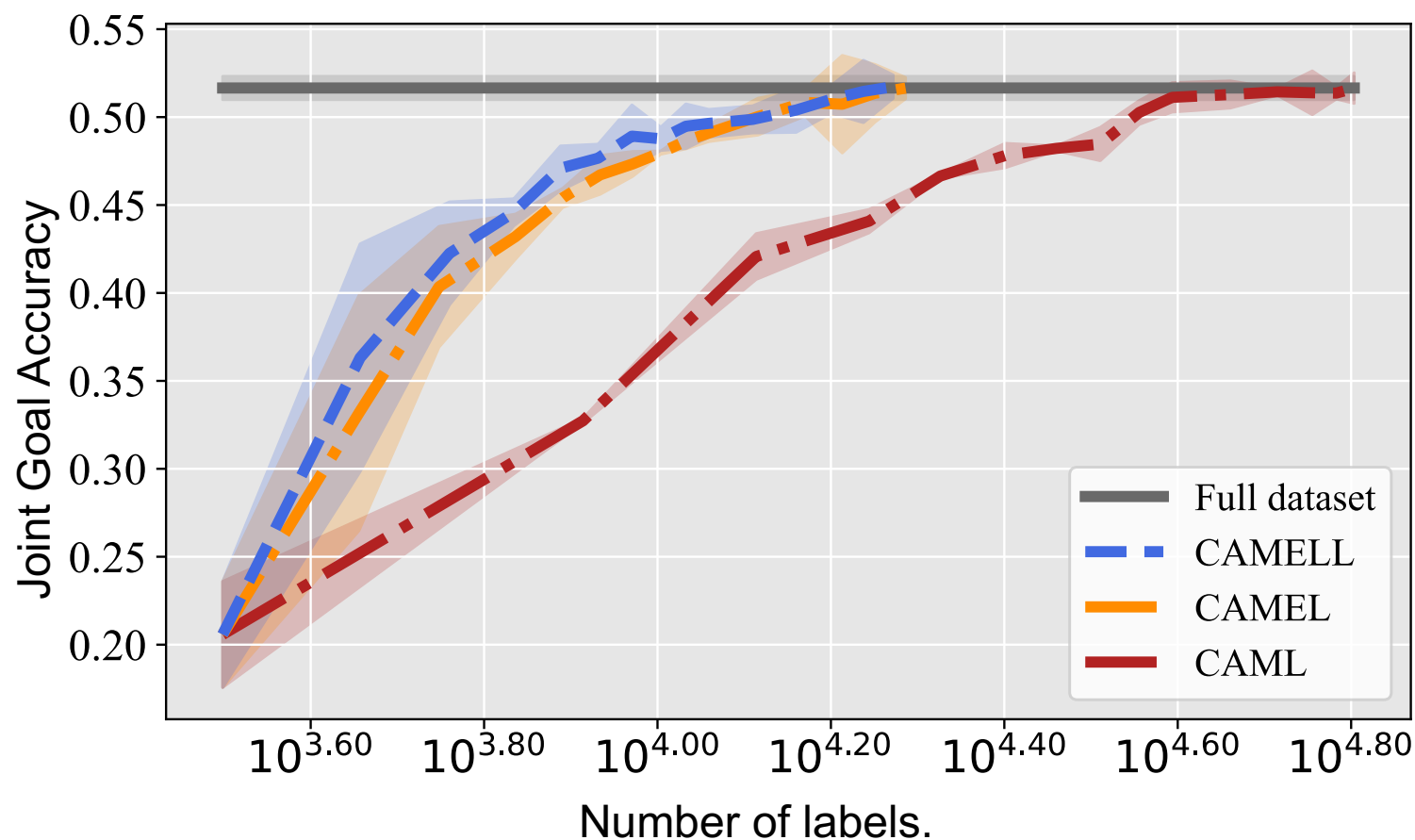


Results

Dialogue Belief Tracking - Ensemble

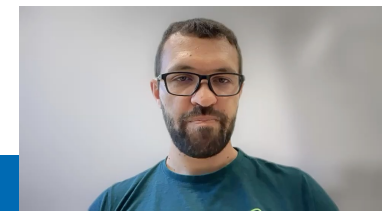
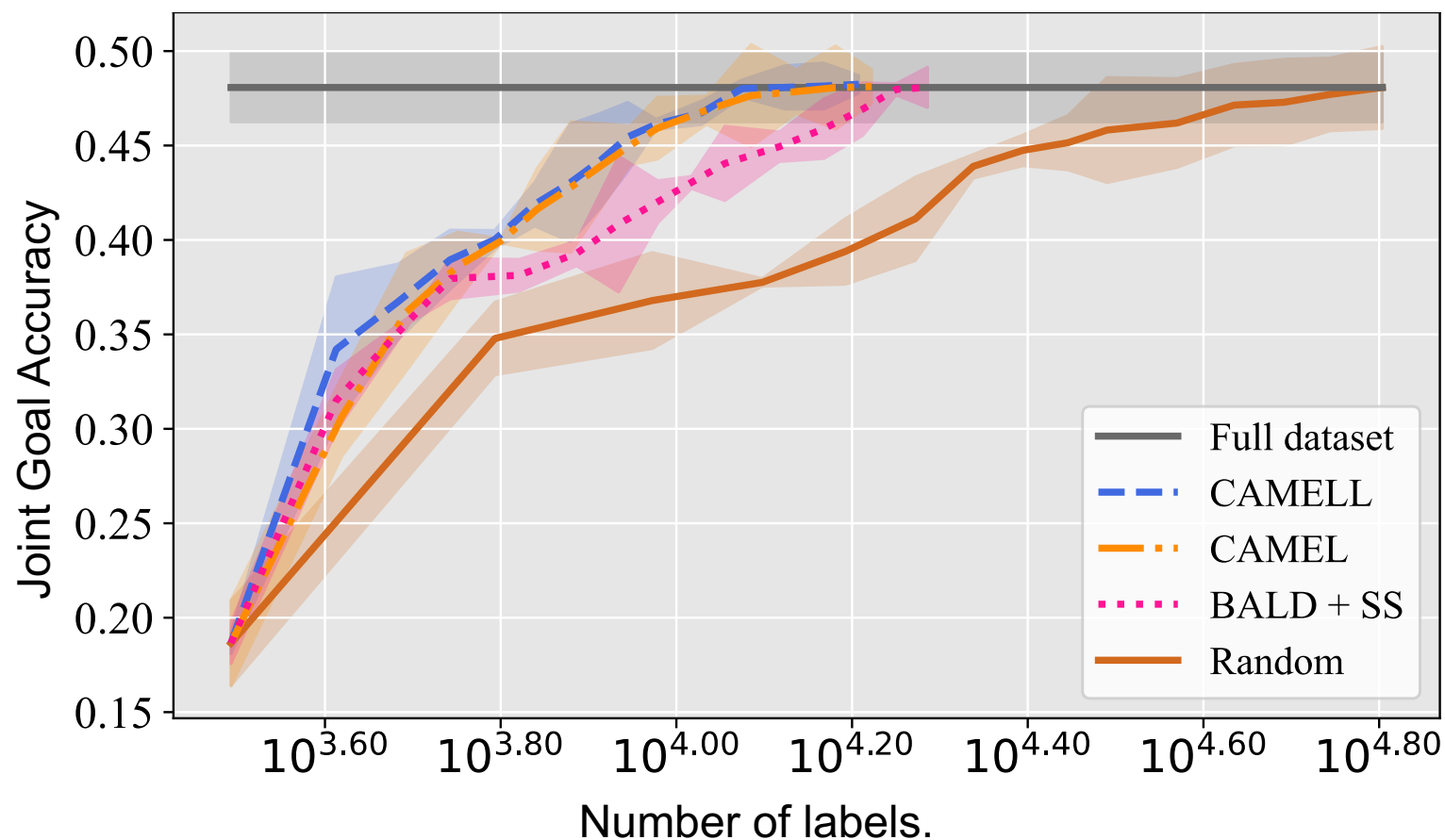


Dialogue Belief Tracking – Ensemble Ablation



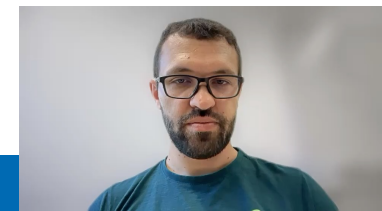
Results

Dialogue Belief Tracking – PostHoc Meta Model

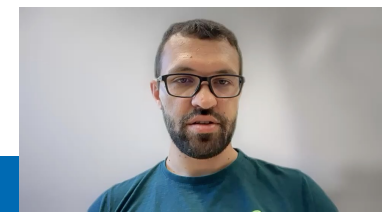


Label Correction Process

- Steps:
 - Select labels in the dataset with label confidence below the threshold.
 - If the prediction confidence is greater than the label threshold replace the label.
- Noisy MultiWOZ 2.1 dataset used to train the ensemble SetSUMBT model.
- Ensemble SetSUMBT used for label correction.
- Trippy (a state-of-the-art span prediction based) DST used to evaluate the quality of the corrections.

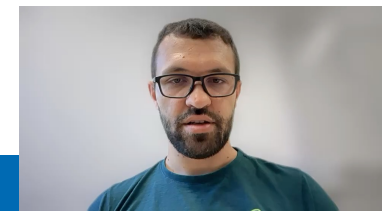


Model	Label Corr.	MultiWOZ 2.1	MultiWOZ 2.4
CE-SetSUMBT	None	51.79	61.63
	Offline	52.83	63.32
TripPy	None	55.28	64.45
	Offline	56.11	66.02



Examples

Conversation	MultiWOZ 2.1 Labels and Corrections
<i>User:</i> I would like to find a place that serves moderately priced Chinese food.	{Restaurant: {Food: Chinese, (95%) Price: Moderate, (94%) Day: Tuesday, (11%) Day: not_mentioned}} (72%)
<i>User:</i> I need a train leaving on Friday and I want to get there by 21:30 . Leaving Broxbourne .	{Train: {Dept.: Broxbourne, (94%) Day: Friday, (95%) Arrive by: 21:20, (1%) Arrive by: 21:30}} (83%)



- ✓ Selective self-supervision **improves efficiency**.
- ✓ CAMELL SetSUMBT achieve 95% of a tracker's full-training dataset performance using **merely 16% of the expert-provided labels**.
- ✓ CAMELL can be applied to **automatically correct labels** in a dataset.

