

# Data-Efficient Drug Design with Pretrained-BERT and Bayesian Active Learning

Muhammad Arslan Masood<sup>1</sup>, Samuel Kaski<sup>1,2</sup>, Tainyu Cui<sup>1,3</sup>,

<sup>1</sup>Department of Computer Science, Aalto University, Finland

<sup>2</sup>University of Manchester, United Kingdom

<sup>3</sup>Imperial College London, United Kingdom

arslan.masood@aalto.fi



## Abstract

In drug discovery, selecting the right compounds to test is critical. Traditional active learning relies only on labeled data, ignoring the rich information in unlabeled molecules. We change that by integrating a BERT model pretrained on 1.26 million unlabeled compounds into the active learning pipeline, enabling robust molecular representations. This significantly improves uncertainty estimation, allowing us to identify toxic compounds with 50% fewer iterations—faster, smarter, and more efficient.

## Research Questions

- ▶ How to perform active learning with limited datasets?
- ▶ How to disentangle representation learning from uncertainty estimation?
- ▶ How to perform target domain aware active learning?

## SS-Active learning

- ▶ Molecular Representations of SMILES  $s_i$  of both labeled  $\mathcal{D}_{\text{labeled}} = \{(s_i^{\text{labeled}}, y_i)\}_{i=1}^N$  and unlabeled set  $\mathcal{D}_U = \{s_i^U\}_{i=1}^{N_U}$  are generated by using pretrained MolBERT model:

$$\mathbf{x}_i^{\text{labeled}} = \text{MolBERT}(s_i^{\text{labeled}}; \Theta_{\text{pretrain}}^*)$$

$$\mathbf{x}_i^U = \text{MolBERT}(s_i^U; \Theta_{\text{pretrain}}^*)$$

$$\mathbf{x}_i^{\text{labeled}}, \mathbf{x}_i^U \in \mathbb{R}^d$$

- ▶ The labelled set  $\mathcal{D}_{\text{labeled}} = \{(\mathbf{x}_i^{\text{labeled}}, y_i)\}_{i=1}^N$  is used to train a probabilistic model  $f(\mathbf{x}_i^{\text{labeled}}; \phi)$  with parameters  $\phi$

$$p(\phi | \mathcal{D}_{\text{labeled}}) \propto p(\mathcal{D}_{\text{labeled}} | \phi) p(\phi)$$

- ▶ A new data point is selected from the unlabeled set  $\mathcal{D}_U = \{(\mathbf{x}_i^U)\}_{i=1}^{N_U}$  by maximizing the acquisition function:

$$\mathbf{x}_s^U = \arg \max_{\mathbf{x} \in \mathcal{D}_U} a(\mathbf{x})$$

- ▶ The label  $y_s$  is acquired for  $\mathbf{x}_s^U$  and incorporated into the training set

$$\mathcal{D}_{\text{updated}} = \mathcal{D} \cup \{(\mathbf{x}_s^U, y_s)\}$$

- ▶ The posterior of probabilistic model is updated by using  $\mathcal{D}_{\text{updated}}$

$$p(\phi | \mathcal{D}_{\text{updated}}) \propto p(\mathcal{D}_{\text{updated}} | \phi) p(\phi)$$

## Acquisition functions

- ▶ Random Acquisition =  $\mathbf{x}_s^U \sim \text{Uniform}(\mathcal{D}_U)$
- ▶ BALD

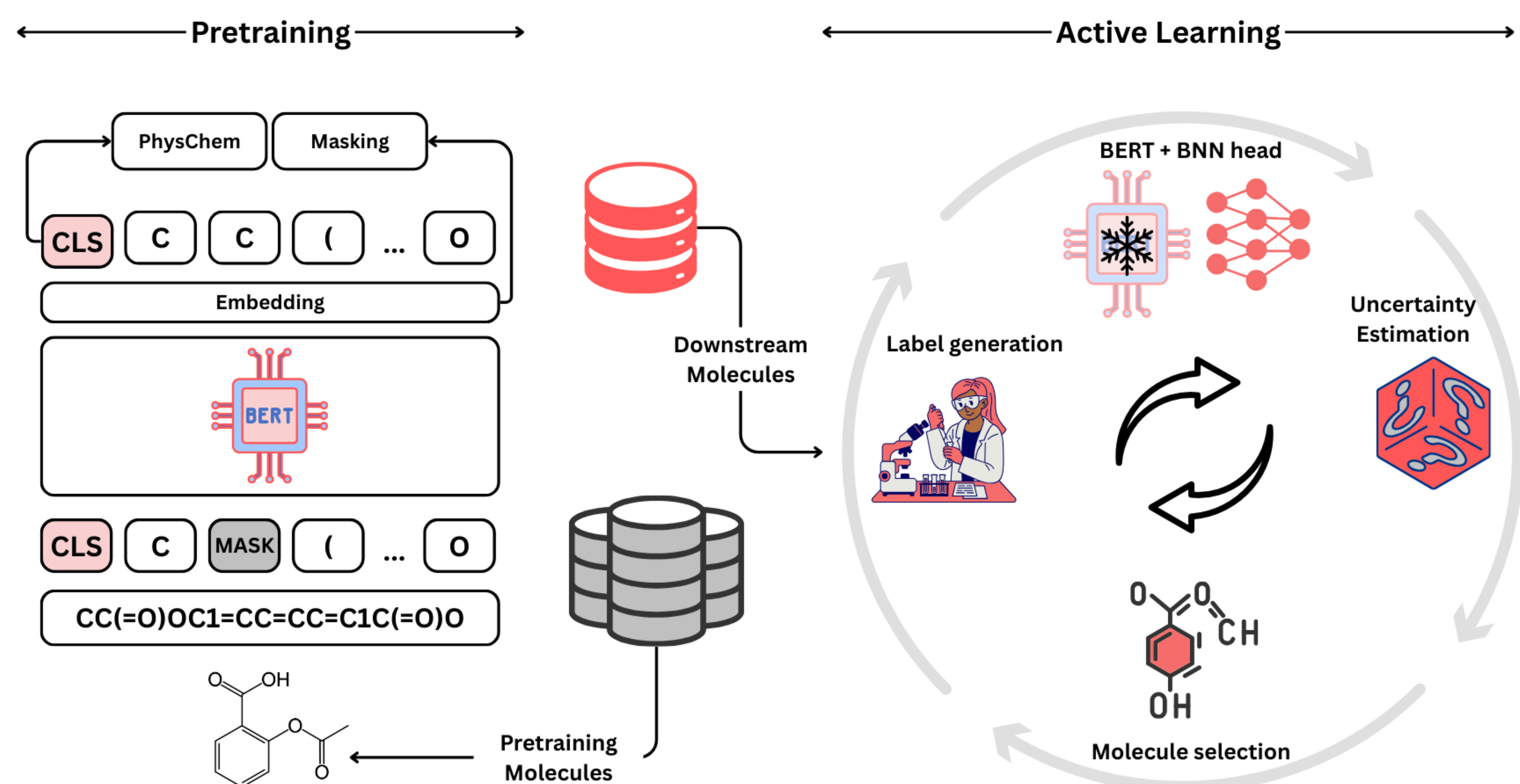
$$= \mathbb{E}_{y \sim p(y|\mathbf{x}, \mathcal{D})} [\mathcal{H}(\phi | \mathcal{D}) - \mathcal{H}(\phi | \mathbf{x}, y, \mathcal{D})]$$

$$= \mathcal{H}[y | \mathbf{x}, \mathcal{D}] - \mathbb{E}_{\phi \sim p(\phi | \mathcal{D})} [\mathcal{H}[y | \mathbf{x}, \phi]]$$

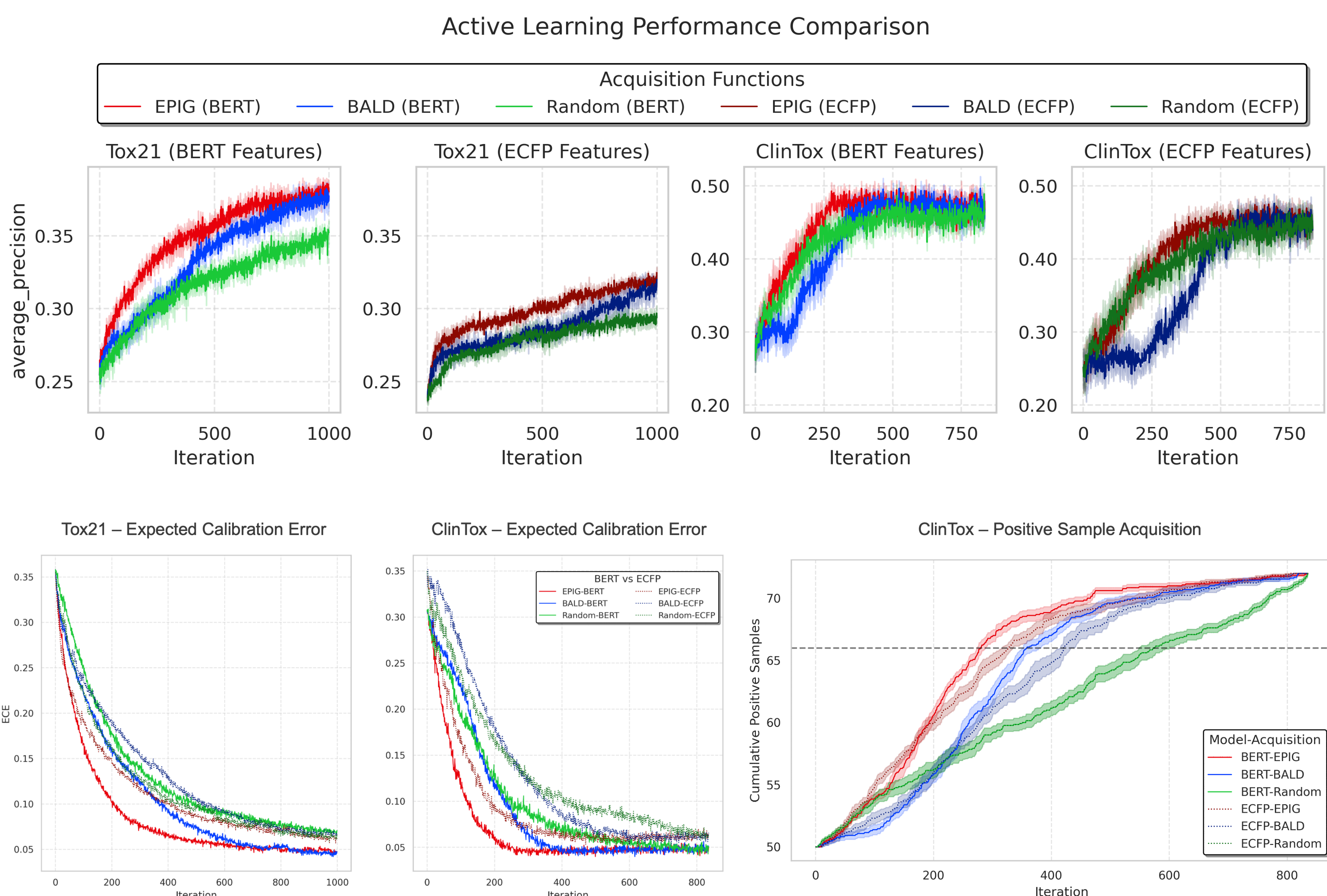
- ▶ EPIG

$$= \mathbb{E}_{p(\mathbf{x}_*)} [\mathcal{H}[y_* | \mathbf{x}_*, \mathcal{D}] - \mathbb{E}_{p(y|\mathbf{x}, \mathcal{D})} [\mathcal{H}[y_* | \mathbf{x}_*, y, \mathbf{x}]]]$$

## Active Learning



## Results



## Acknowledgements

This project is funded by Horizon2020 research and innovation programme under the MSC-ITN actions, grant agreement No. 956832.