

Automatic Domain Adaptation by Transformers in In-Context Learning

<https://arxiv.org/abs/2405.16819>

Ryuichiro Hataya ^{α} , Kota Matsui ^{β} , Masaaki Imaizumi ^{γ} ^{α}

^{α} RIKEN AIP, ^{β} Nagoya University, ^{γ} The University of Tokyo

contact: ryuichiro.hataya@riken.jp / webpage: hataya.tokyo

Summary

- Selecting or designing suitable UDA algorithms for given problems is challenging
- We showed that Transformers can
 - approximate UDA algorithms (IWL and DANN)
 - select appropriate one based on data statistics in the ICL framework

Unsupervised Domain Adaptation (UDA)

Source

labeled data



$$\mathcal{D}_S = \{(x_i^S, y_i^S)\}_{i=1}^n$$
$$(x_i^S, y_i^S) \sim p_S$$

Target

unlabeled data



$$\mathcal{D}_T = \{x_i^T\}_{i=1}^{n'}$$
$$(x_i^T, \cdot) \sim p_T$$

transfer “knowledge” of source
to classify target data

Goal: minimize the target risk

i.e., get $\operatorname{argmin}_{f \in \mathcal{F}} R(f)$, where $R(f) = \mathbb{E}_{(x, y) \sim p_T} [\ell(f(x), y)]$

Unavailable!

UDA: Instance-based Approach

When $\mathcal{X}_S \times \mathcal{Y}_S = \mathcal{X}_T \times \mathcal{Y}_T$ and $p_S(\mathbf{x}, y) \neq p_T(\mathbf{x}, y)$

Assume **covariate shift** $p_S(y|\mathbf{x}) = p_T(y|\mathbf{x})$

Importance-weighted learning [Sugiyama+12, Kimura&Hino24]

$$R(f) = \mathbb{E}_{(\mathbf{x}, y) \sim p_T} [\ell(f(\mathbf{x}, y))] = \mathbb{E}_{(\mathbf{x}, y) \sim p_S} [q(\mathbf{x}) \ell(f(\mathbf{x}), y)] , \text{ where } q(\mathbf{x}) = \frac{p_T(\mathbf{x})}{p_S(\mathbf{x})}$$

1. Estimate density ratio by e.g., uLSIF [Kanamori+09]

$$\hat{q}_{\hat{\alpha}} = \hat{\alpha}^\top \phi(\mathbf{x}) , \text{ where } \hat{\alpha} = \operatorname{argmin}_{\alpha} \frac{1}{2} \int (\alpha^\top \phi(\mathbf{x}) - q(\mathbf{x}))^2 p_S(\mathbf{x}) d\mathbf{x}$$

2. Minimize the target risk

$$\hat{R}(w) = \sum_{(\mathbf{x}, y) \in \mathcal{D}_S} \hat{q}_{\hat{\alpha}}(\mathbf{x}) \ell(w^\top \phi(\mathbf{x}), y)$$

UDA: Feature-based Approach

When $\mathcal{X}_S \times \mathcal{Y}_S \neq \mathcal{X}_T \times \mathcal{Y}_T$

Find domain invariant features ψ s.t. $R(f') \approx \mathbb{E}_{(x,y) \sim p_S} [\ell(f'(\psi(x)), y)]$

Domain Adversarial Neural Networks (DANNs) [Ganin+16]

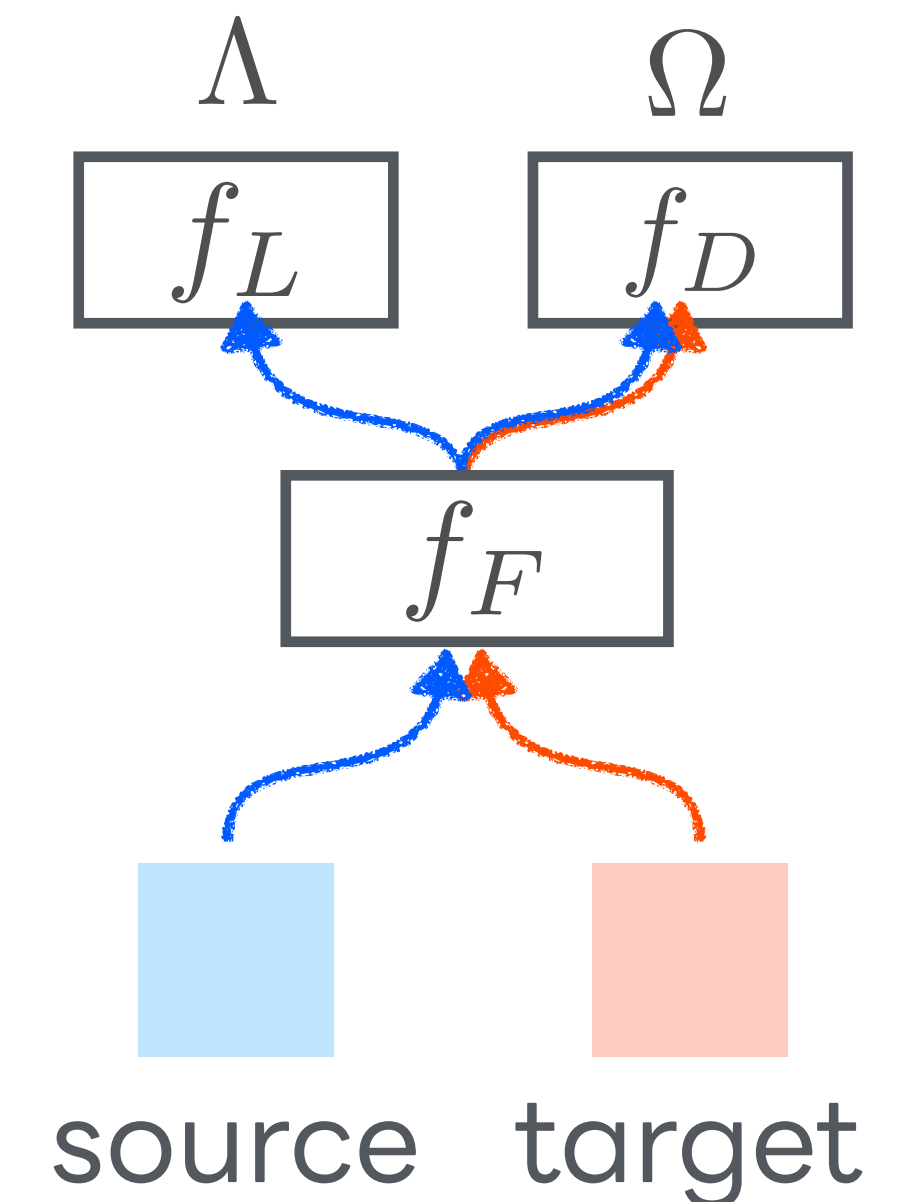
f_F : feature extractor, parameterized by θ_F (ψ)

f_L : label classifier, parameterized by θ_L (f')

f_D : domain classifier, parameterized by θ_D

DANN obtains ψ and f' by adversarial learning

$$\min_{\theta_F, \theta_L} \max_{\theta_D} \Lambda(\theta_F, \theta_L) - \lambda \Omega(\theta_F, \theta_D)$$



In-context Learning (ICL)

LLMs perform new tasks from instructions in prompts



cat: 0
dog: 0
whale: 1
dolphin: 1
cod: 0
salmon: 1
eel: ?

} training data

} test data

→ LLMs perform “machine learning”
without parameter update



ChatGPT

0

ICL: Theory

Appropriately pre-trained TFs can perform ML algorithms

- Least squares [Zhang+23, Akyürek+23]
- Gradient descent of linear regression [von Oswald+22, Akyürek+23]

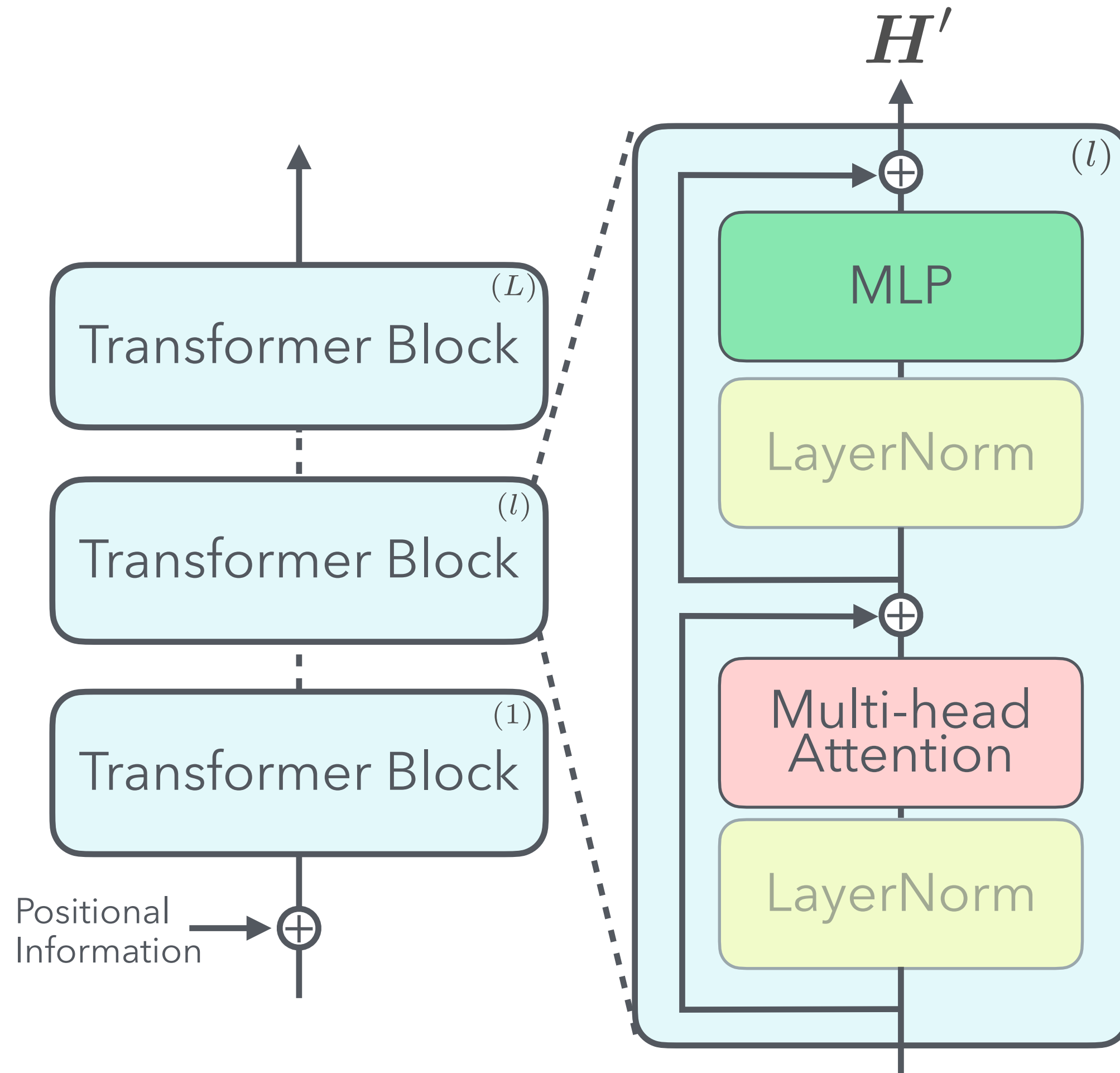
GLMs / neural networks [Bai+23]

- Linear UCB, Thompson sampling, ... [Lin+24]

TFs can also perform model selection [Bai+23]

Based on validation performance / data types

Transformer



$$\text{MLP}^{(l)}(\mathbf{H}) = \mathbf{H} + \mathbf{W}_2^{(l)} \varsigma(\mathbf{W}_1^{(l)} \mathbf{H})$$

intra-token transformation

$\mathbf{W}_1^{(l)} \in \mathbb{R}^{D' \times D}$, $\mathbf{W}_2^{(l)} \in \mathbb{R}^{D \times D'}$: learnable parameters

$$\text{Attn}^{(l)}(\mathbf{H}) = \mathbf{H} + \frac{1}{N} \sum_{m=1}^M \mathbf{V}_m^{(l)} \mathbf{H} \sigma((\mathbf{Q}_m^{(l)} \mathbf{H})^\top \mathbf{K}_m^{(l)} \mathbf{H})$$

inter-token transformation

$\mathbf{K}_m^{(l)}, \mathbf{Q}_m^{(l)}, \mathbf{V}_m^{(l)} \in \mathbb{R}^{D \times D}$: learnable parameters

$$\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_n, \dots, \mathbf{h}_N] \in \mathbb{R}^{N \times D}$$

token

ς σ : activation function (ReLU in this work)

ICL Framework by Bai+23

1. Sum-of-ReLUs can approximate any functions

e.g., $f(x, y) = \sum_{m=1}^M c_m \sigma(a_m x + b_m y + d_m)$ can approximate any bivariate functions

2. Multi-head Attention can represent any aggregation functions

$$\text{Attn}^{(l)}(\mathbf{h}_i) = \mathbf{h}_i + \frac{1}{N} \sum_{j=1}^N \sum_{m=1}^M \sigma(\langle \mathbf{Q}_m^{(l)} \mathbf{h}_i, \mathbf{K}_m^{(l)} \mathbf{h}_j \rangle) \mathbf{V}_m^{(l)} \mathbf{h}_j$$

aggregation
e.g., mean of selected data

approximate a function

j s with positive values in σ are selected

3. Single attention block can imitate a single GD step of GLM

Let $\mathbf{h}_i = \begin{bmatrix} \mathbf{x}_i \\ y_i \\ \mathbf{w} \end{bmatrix}$

$$\text{Attn}^{(l)}(\mathbf{h}_i) = \begin{bmatrix} \mathbf{x}_i \\ y_i \\ \mathbf{w} \end{bmatrix} + \frac{1}{N} \sum_{j=1}^N \sum_{m=1}^M \sigma(\langle \mathbf{Q}_m^{(l)} \mathbf{h}_i, \mathbf{K}_m^{(l)} \mathbf{h}_j \rangle) \mathbf{V}_m^{(l)} \mathbf{h}_j$$

$a_m \mathbf{w}^\top \mathbf{x}_j + b_m y_j + d_m$

$-\eta c_m \begin{bmatrix} 0 \\ 0 \\ \mathbf{x}_j \end{bmatrix}$

$\approx \begin{bmatrix} 0 \\ 0 \\ -\eta \nabla_{\mathbf{w}} \hat{L} \end{bmatrix}$

$\approx -\eta \partial_1 \ell(\mathbf{w}^\top \mathbf{x}_j, y_j) \begin{bmatrix} 0 \\ 0 \\ \mathbf{x}_j \end{bmatrix}$

$$\hat{L} = \sum_{j=1}^N \ell(\mathbf{w}^\top \mathbf{x}_j, y_j)$$

model parameter

loss function

ICL Framework by Bai+23

3. Single attention block can imitate a single GD step of GLM

Let $\mathbf{h}_i = \begin{bmatrix} \mathbf{x}_i \\ y_i \\ \mathbf{w} \end{bmatrix}$

$$\text{Attn}^{(l)}(\mathbf{h}_i) = \begin{bmatrix} \mathbf{x}_i \\ y_i \\ \mathbf{w} \end{bmatrix} + \frac{1}{N} \sum_{j=1}^N \sum_{m=1}^M \sigma(\langle \mathbf{Q}_m^{(l)} \mathbf{h}_i, \mathbf{K}_m^{(l)} \mathbf{h}_j \rangle) \mathbf{V}_m^{(l)} \mathbf{h}_j$$

Diagram illustrating the attention block structure and its relationship to the GD step of GLM:

- The input \mathbf{h}_i is a vector containing \mathbf{x}_i , y_i , and \mathbf{w} .
- The attention block computes a weighted sum over N samples and M heads.
- The weight is $\sigma(\langle \mathbf{Q}_m^{(l)} \mathbf{h}_i, \mathbf{K}_m^{(l)} \mathbf{h}_j \rangle)$.
- The value is $\mathbf{V}_m^{(l)} \mathbf{h}_j$.
- The output is $\begin{bmatrix} \mathbf{x}_i \\ y_i \\ \mathbf{w} \end{bmatrix} + \frac{1}{N} \sum_{j=1}^N \sum_{m=1}^M \sigma(\langle \mathbf{Q}_m^{(l)} \mathbf{h}_i, \mathbf{K}_m^{(l)} \mathbf{h}_j \rangle) \mathbf{V}_m^{(l)} \mathbf{h}_j$.
- The diagram shows the relationship between the attention block and the GD step of GLM:

Let $\hat{L} = \sum_{j=1}^N \ell(\mathbf{w}^\top \mathbf{x}_j, y_j)$

Annotations:

- $\mathbf{w}^\top \mathbf{x}_j$: model parameter
- $\ell(\mathbf{w}^\top \mathbf{x}_j, y_j)$: loss function
- $\mathbf{Q}_m^{(l)} \mathbf{h}_i$: $\approx \begin{bmatrix} 0 \\ 0 \\ -\eta \nabla_{\mathbf{w}} \hat{L} \end{bmatrix}$
- $\mathbf{K}_m^{(l)} \mathbf{h}_j$: $\approx -\eta \partial_1 \ell(\mathbf{w}^\top \mathbf{x}_j, y_j) \begin{bmatrix} 0 \\ 0 \\ \mathbf{x}_j \end{bmatrix}$
- $\mathbf{V}_m^{(l)} \mathbf{h}_j$: $\approx \begin{bmatrix} 0 \\ 0 \\ \mathbf{x}_j \end{bmatrix}$

4. L -layer TF can approximate L steps of GD

Let the input be $\mathbf{H}^{(1)} =$

\mathbf{x}_1	\mathbf{x}_2	\dots	\mathbf{x}_n	\mathbf{x}_*
y_1	y_2	\dots	y_n	0
1	1	\dots	1	1
1	1	\dots	1	0
0	0	\dots	0	0

Marking training data

training data test data

For any $\varepsilon > 0$, there exists an L -layer TF s.t. $\|\text{read}(\text{TF}(\mathbf{H}^{(1)})) - (\mathbf{w}^*)^\top \mathbf{x}_*\| \leq \varepsilon$

model parameter after L steps of GD

Setup: In-context Learning Domain Adaptation

We encode data as

$$H^{(1)} = \begin{bmatrix} \boxed{x_1^S \quad \dots \quad x_n^S} & \boxed{x_1^T \quad \dots \quad x_{n'}^T} & \boxed{x_*} \\ y_1 \quad \dots \quad y_n & 0 \quad \dots \quad 0 & 0 \\ \textcolor{blue}{1} \quad \dots \quad \textcolor{blue}{1} & \textcolor{blue}{0} \quad \dots \quad \textcolor{blue}{0} & \textcolor{blue}{0} \\ 1 \quad \dots \quad 1 & 1 \quad \dots \quad 1 & 1 \\ \textcolor{orange}{1} \quad \dots \quad \textcolor{orange}{1} & \textcolor{orange}{1} \quad \dots \quad \textcolor{orange}{1} & \textcolor{orange}{1} \\ \mathbf{0} \quad \dots \quad \mathbf{0} & \mathbf{0} \quad \dots \quad \mathbf{0} & \mathbf{0} \end{bmatrix}$$

training data from source
training data from target
test data from target

Marking data from source: d_i
 Marking training data

Note: inputs like $[x_1^S, y_1, x_2^S, y_2, \dots, x_1^T, \dots, x_*]$
 can be converted into $H^{(1)}$ using some TF layers

Main results: ICL Domain Adaptation Algorithms

Theorem 1

For any $\varepsilon > 0$, there exists a $2L + 1$ -layer Transformer that can approximate importance-weighted learning with the uLSIF estimator:

$$\|\text{read}(\text{TF}(\mathbf{H}^{(1)})) - \hat{f}^{\text{IWL}}(\mathbf{x}_*)\| \leq \varepsilon$$

Theorem 2

For any $\varepsilon > 0$, there exists a $2L$ -layer Transformer that can approximate a two-layer domain adversarial neural network:

$$\|\text{read}(\text{TF}(\mathbf{H}^{(1)})) - \hat{f}^{\text{DANN}}(\mathbf{x}_*)\| \leq L\varepsilon$$

Transformers can also approximate other UDA algorithms

Main results: ICL Algorithm Selection

Select the “appropriate” result for given data

🤔 No labels for target data. How to select one?

Use IWL if supports of p_S and p_T overlap sufficiently,
otherwise DANN

$$\hat{f}^{\text{ICUDA}}(\mathbf{x}_*) = \begin{cases} \hat{f}^{\text{IWL}}(\mathbf{x}_*) & \text{if } \min_{\mathbf{x} \sim \mathcal{D}_T} p_S(\mathbf{x}) > \delta \\ \hat{f}^{\text{DANN}}(\mathbf{x}_*) & \text{otherwise} \end{cases}$$

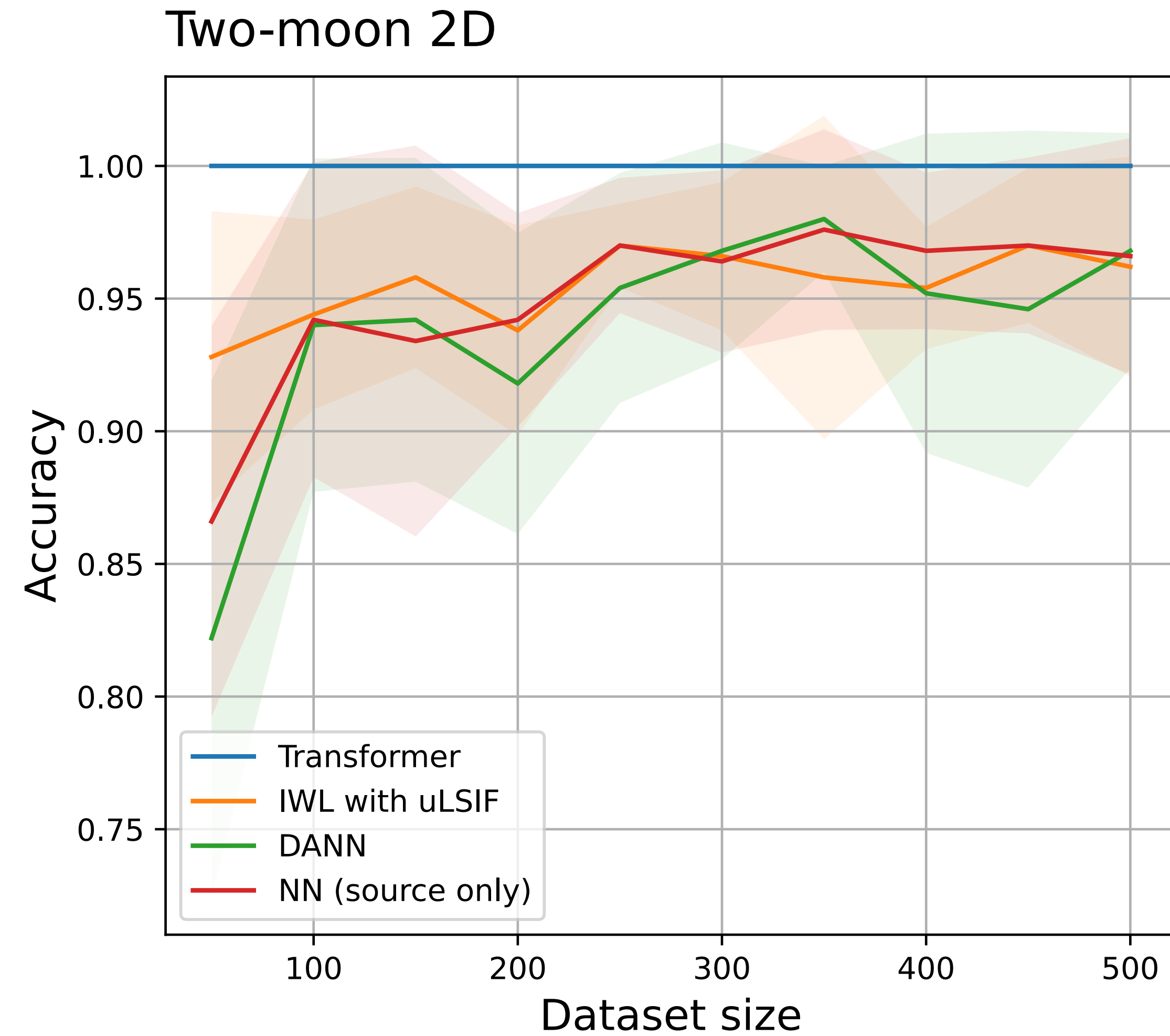
Theorem 3

For any $\varepsilon > 0$, there exists a three-layer Transformer that approximately select the appropriate result with prob. at least $1 - \frac{1}{n'} - O(\varepsilon)$:

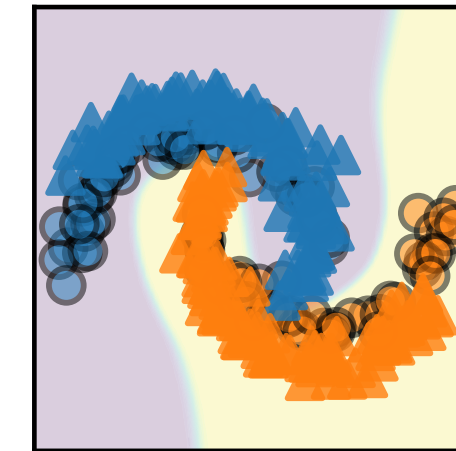
$$\|\text{read}(\text{TF}(\mathbf{H}')) - \hat{f}^{\text{ICUDA}}(\mathbf{x}_*)\| \leq \varepsilon$$

Transformers can also approximate other selection rules

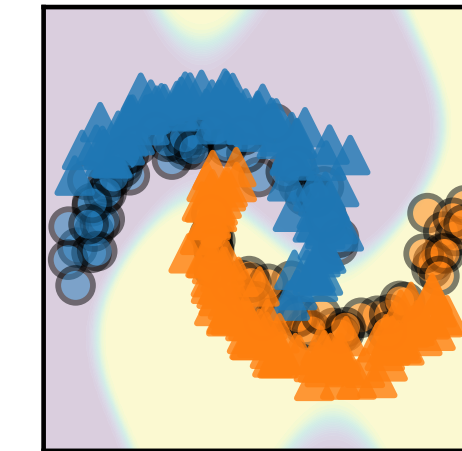
Experiments



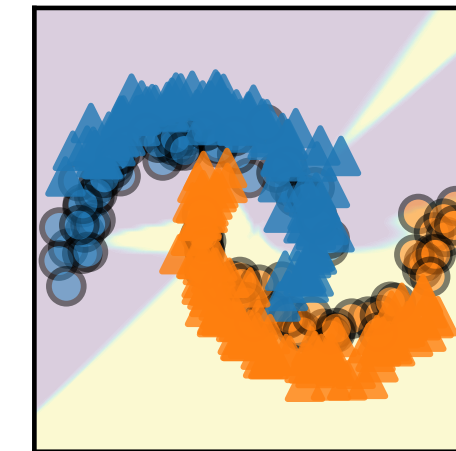
Transformer



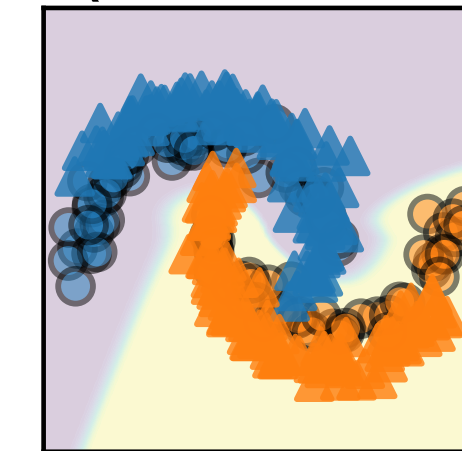
IWL with uLSIF



DANN



NN (source only)



Source Target

Summary

- We showed that Transformers can
 - approximate UDA algorithms (IWL and DANN)
 - select appropriate one based on data statisticsin the ICL framework
- TFs can approximate other algorithms / selection rules
 - Mixing / combining algorithms is also possible→ LLMs may know better UDA algorithms unknown to us