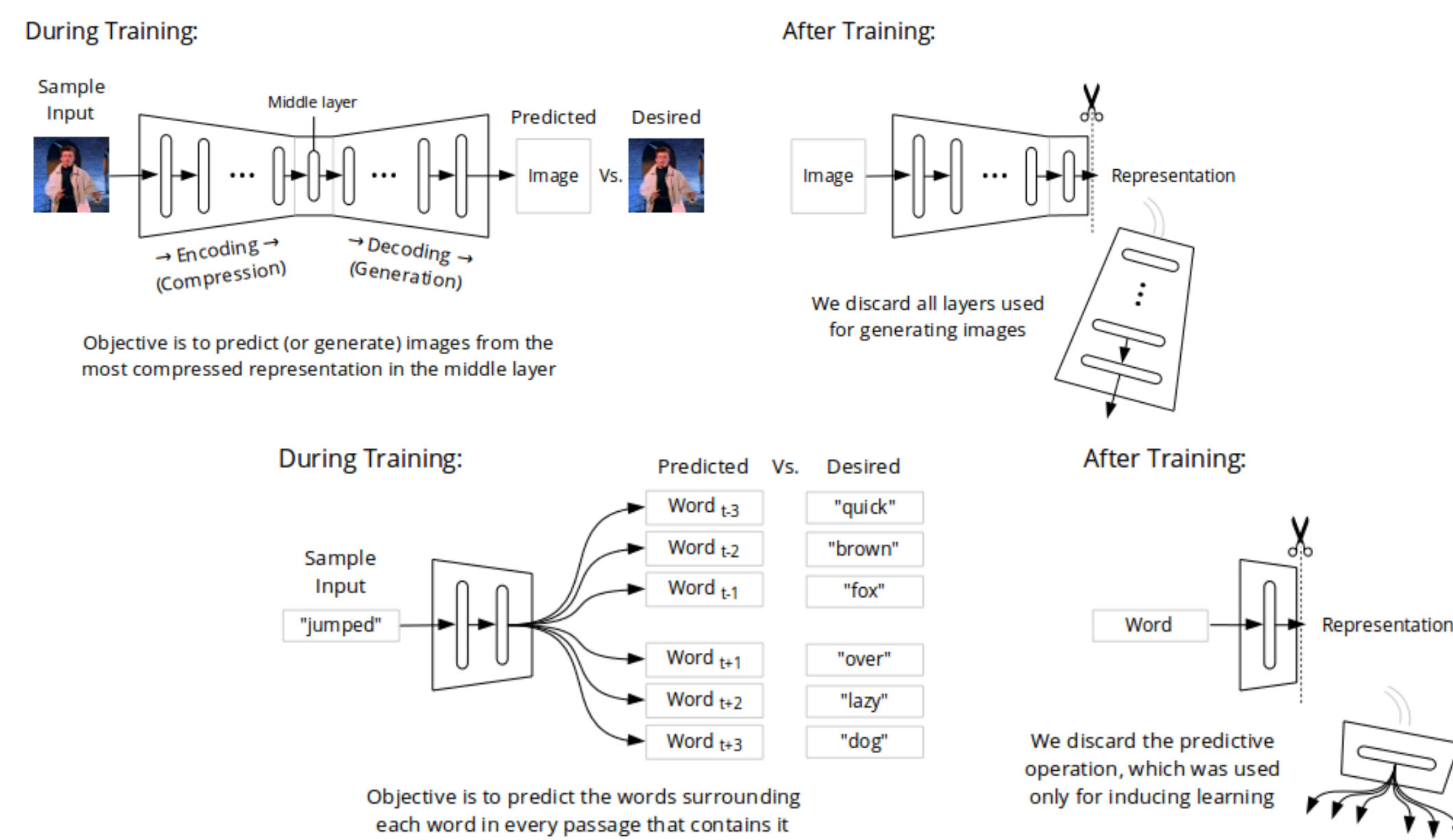
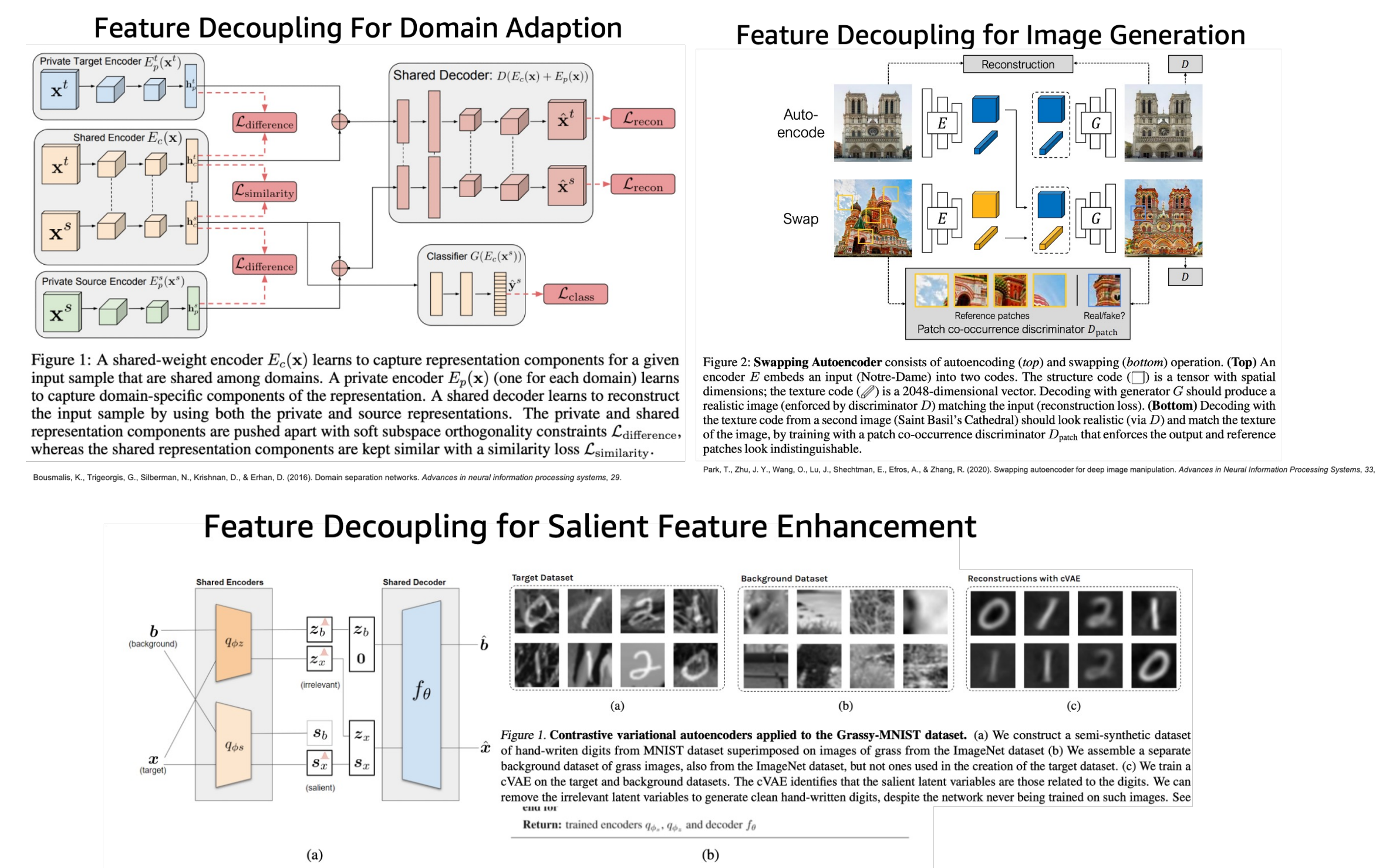


Background

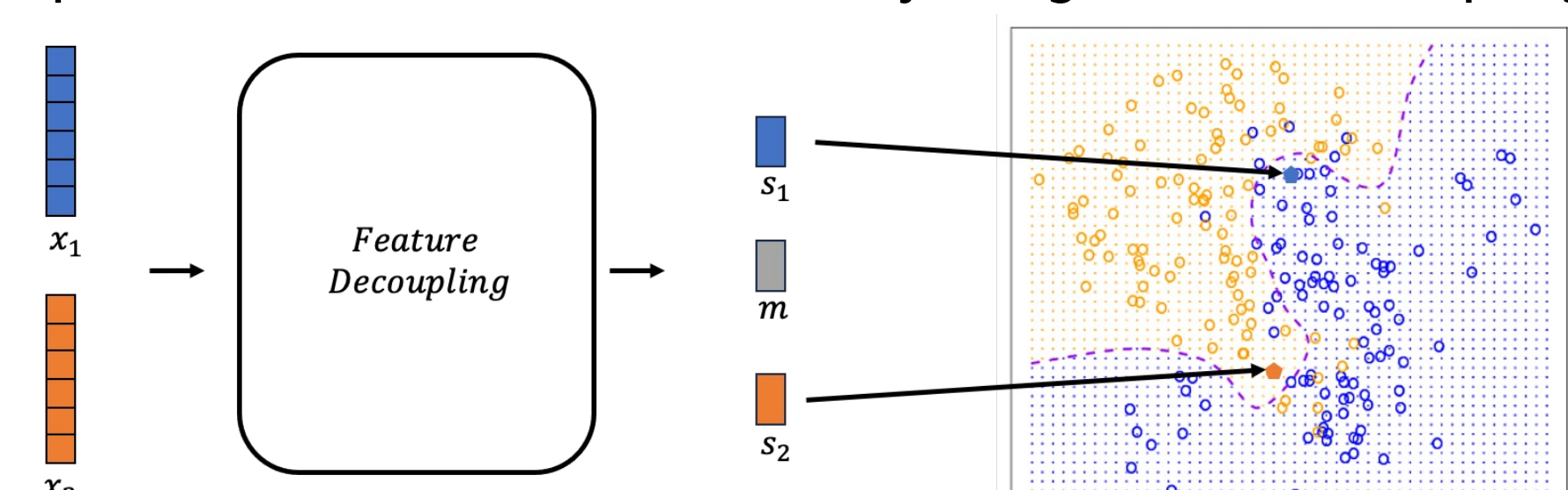
- Representation Learning is to learn to transform data from its original representation to a new representation that retains information essential to objects that are of interest to our tasks
- Representation Learning has made remarkable advancements in Computer Vision and Natural Language Processing



Feature Decoupling



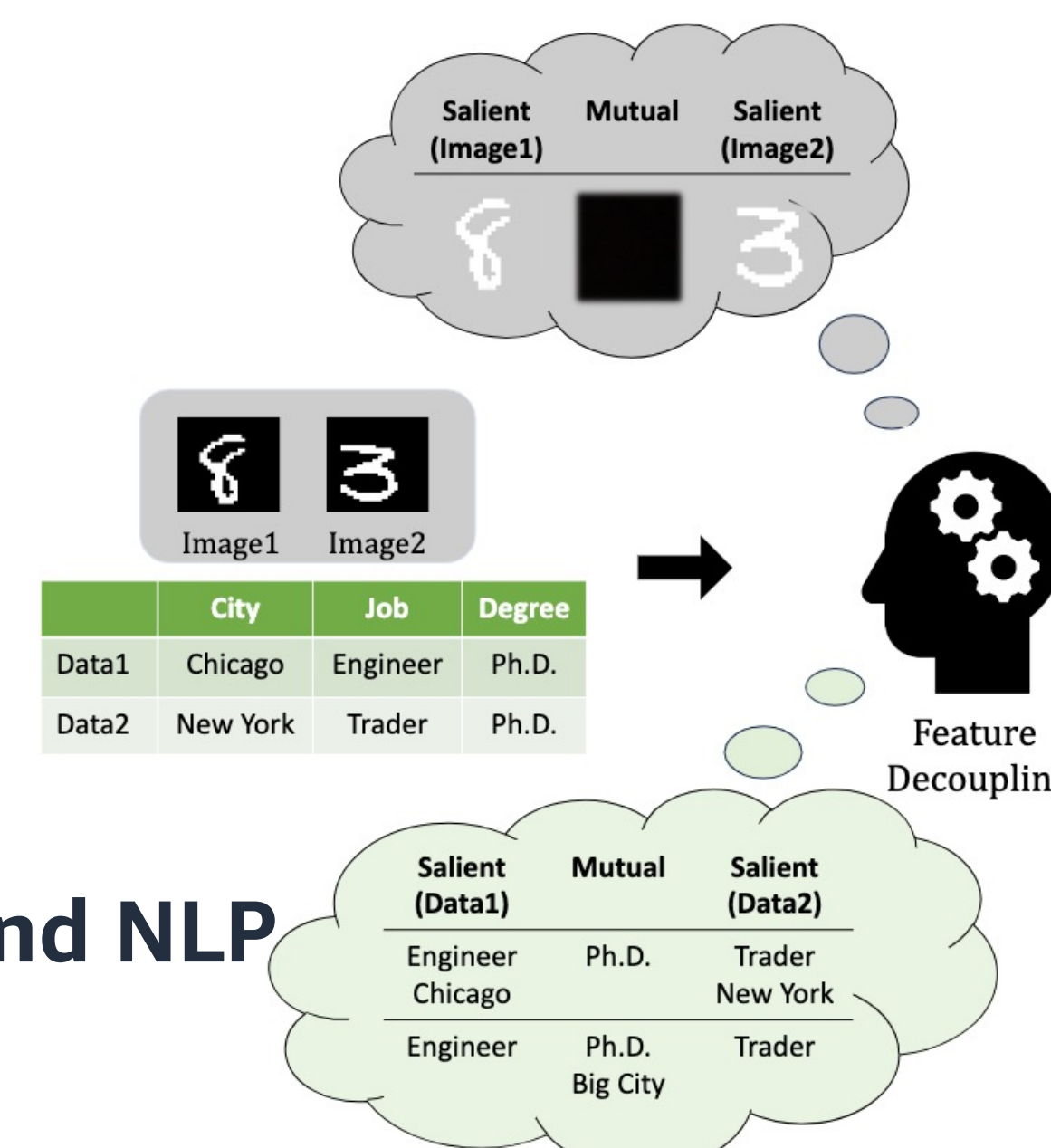
Can we help tabular data to draw a boundary using Feature Decoupling?



- When comparing two data samples,
- mutual features consist of information that highlights common characteristics.
 - salient features emphasize the distinctive attributes to differentiate one sample from the other.

Motivation

- Tabular Representation Learning has not been fully explored, due to
 - Inherent heterogeneity which lacks explicit spatial relationships found in images (e.g., similar background and distinct characters) or the semantic dependencies observed in languages
 - Various discrete and continuous distributions from both numerical and categorical features
 - Complex interrelationship from features that can be dependent or independent from each other



Can we adopt the success of Representation Learning from CV and NLP domains to Tabular data?

- For an image, a person can easily distinguish the salient digits from the mutual background
- Separating the salient and mutual information becomes challenging for tabular samples
- Explicitly distinguish between mutual information in global and salient information within the feature space can be useful for representation learning

Figure 1: Given a pair of images, a person can easily distinguish the salient digits and mutual background due to the well-structured spatial relationships. However, it becomes challenging to distinguish a pair of tabular samples. For instance, feature City may be salient between data points "Chicago" and "New York" for word counts, however, still sharing some latent mutual information (e.g., big cities), making it challenging for decoupling. Note that this decoupling process is for illustration only. In the implementation, all the decoupled samples are computed in the feature space.

SwitchTab Self-supervised Learning Framework

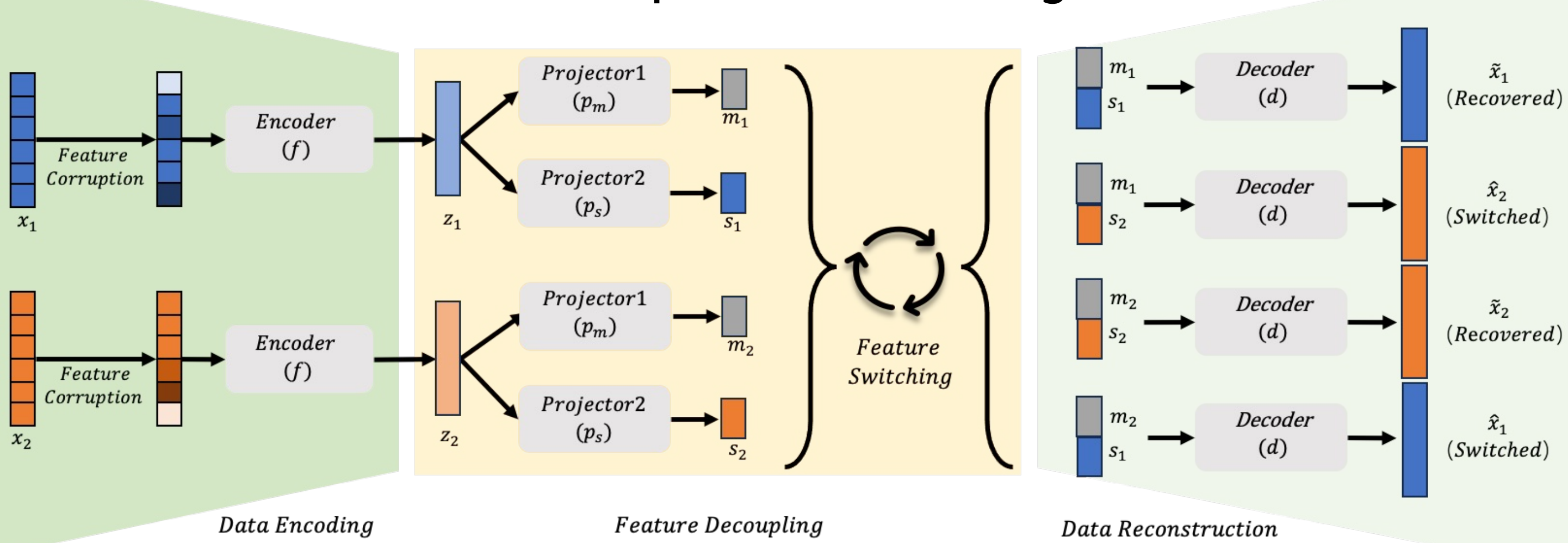
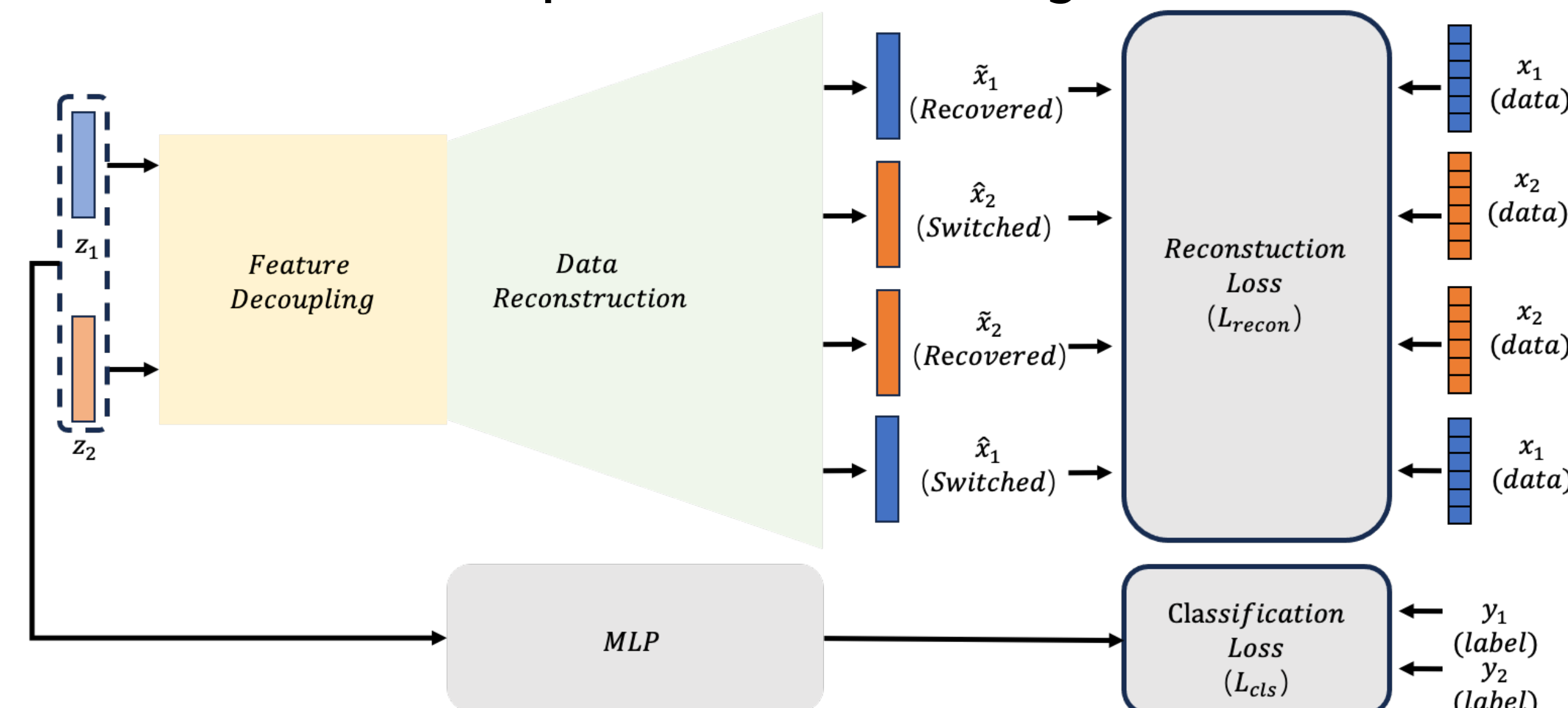


Figure 2: Block diagram of the proposed self-supervised learning framework. (1) Two different samples x_1 and x_2 are randomly corrupted and encoded into feature vectors z_1 and z_2 through encoder f . (2) feature vectors z_1 and z_2 are decoded into mutual and salient features by two different projectors p_m and p_s , respectively. (3) Mutual and salient features are combined and reconstructed by a decoder d where the salient feature dominates the sample type and the mutual feature provides common information that is switchable among two samples.

SwitchTab Supervised Pretraining with Labels



Loss Function and Objectives

Algorithm 1: Self-supervised Learning with SwitchTab

Require: unlabeled data $\mathcal{X} \subseteq \mathbb{R}^{d \times d}$, batch size B , encoder f , projector for mutual information p_m , projector for salient information p_s , decoder d , mean squared error MSE, feature concatenation \oplus .

- for two sampled mini-batch $\{x_i^a\}_{i=1}^B \subseteq \mathcal{X}$ and $\{x_j^b\}_{j=1}^B \subseteq \mathcal{X}$ do
- for each sample x_i^a and x_j^b , apply feature corruption, define the corrupted feature as \tilde{x}_i^a and \tilde{x}_j^b , for $i \in [B]$
- data encoding: $z_i^a = f(x_i^a), z_j^b = f(x_j^b)$, for $i \in [B]$
- feature decoupling:
 - the salient and mutual information of the first batch be defined as follows: $s_i^a = p_s(z_i^a)$ and $m_i^a = p_m(z_i^a)$.
 - the salient and mutual information of the second batch be defined as follows: $s_j^b = p_s(z_j^b)$ and $m_j^b = p_m(z_j^b)$.
- data reconstruction:
 - let recovered pairs be defined as: $\tilde{x}_i^a = d(m_i^a \oplus s_i^a), \tilde{x}_j^b = d(m_j^b \oplus s_j^b)$
 - let switched pairs be defined as: $\tilde{x}_i^a = d(m_j^b \oplus s_i^a), \tilde{x}_j^b = d(m_i^a \oplus s_j^b)$
- define reconstruction loss $L_{recon} = \text{MSE}(x_i^a, \tilde{x}_i^a) + \text{MSE}(x_j^b, \tilde{x}_j^b) + \text{MSE}(x_i^a, \tilde{x}_j^b) + \text{MSE}(x_j^b, \tilde{x}_i^a)$
- update encoder f , projectors p_m and p_s , and decoder d to minimize L_{recon} using RMSProp.
- end for

$$L_{recon} = \frac{1}{M} \sum_{j=1}^M (x_{1j} - \tilde{x}_{1j})^2 + \frac{1}{M} \sum_{j=1}^M (x_{2j} - \tilde{x}_{2j})^2 + \frac{1}{M} \sum_{j=1}^M (x_{1j} - \tilde{x}_{2j})^2 + \frac{1}{M} \sum_{j=1}^M (x_{2j} - \tilde{x}_{1j})^2$$

$$L_{cls} = L_{recon} + \alpha \cdot L_{cls}$$

where α is used to balance the classification loss and reconstruction loss and set to 1 as default. To illustrate, the cross-entropy loss used for classification task can be defined as follows:

$$L_{cls} = - (y_1 \log(\hat{y}_1) + y_2 \log(\hat{y}_2))$$

where \hat{y}_1 and \hat{y}_2 are predicted labels computing a MLP, i.e., $\hat{y}_1 = \text{MLP}(z_1)$ and $\hat{y}_2 = \text{MLP}(z_2)$. For regression tasks, we replace the cross-entropy loss with root mean squared error (RMSE).

Case Study Experiments

- Dataset
- A standard benchmark from (Gorishniy et al. 2021) with 11 datasets including California Housing (CA), Adult (AD), Helena (HE), Jannis (JA), Higgs (HI), ALOI (AL), Epsilon (EP), Year (YE), Covertypes (CO), Yahoo (YA), Microsoft(MI)
 - Additional popular datasets from recent work with 7 datasets on classification tasks
 - Bank (BK), Blastchar (BC), Arrhythmia (AT), Arcene (AR), Shoppers (SH), Volkert (VO), MNIST (MN)

Dataset size	48842	65196	83733	98050	108000	500000	518012	20640	515345	709877	1200192
Feature size	14	27	54	28	128	2000	54	90	12	699	136
Method/Dataset	AD ↑	HE ↑	JA ↑	HI ↑	AL ↑	EP ↑	CO ↑	CA ↓	YA ↓	MI ↓	
TabNet	0.850	0.378	0.723	0.719	0.954	0.8896	0.957	0.510	8.909	0.823	0.751
SNN	0.854	0.373	0.719	0.722	0.954	0.8975	0.961	0.493	8.895	0.761	0.751
AutoInt	0.859	0.372	0.721	0.725	0.945	0.8949	0.934	0.474	8.882	0.768	0.750
MLP	0.852	0.363	0.723	0.723	0.954	0.8977	0.962	0.499	8.853	0.757	0.747
DCN2	0.853	0.385	0.723	0.723	0.955	0.8977	0.965	0.484	8.800	0.757	0.749
NODE	0.858	0.359	0.726	0.726	0.918	0.8958	0.985	0.464	8.784	0.753	0.745
ResNet	0.854	0.396	0.727	0.727	0.963	0.8969	0.964	0.486	8.846	0.757	0.748
CatBoost	0.873	0.388	0.727	0.729	0.948	0.8993	0.950	0.423	8.837	0.740	0.743
FT-Transformer	0.859	0.391	0.729	0.729	0.960	0.8982	0.970	0.459	8.855	0.756	0.746
XGBoost	0.874	0.377	0.724	0.728	0.924	0.8799	0.964	0.431	8.819	0.732	0.742
SwitchTab (Self-Sup.)	0.867	0.387	0.726	0.724	0.942	0.8928	0.971	0.452	8.857	0.755	0.751
SwitchTab	0.881	0.389	0.731	0.733	0.951	0.8987	0.989	0.442	8.822	0.744	0.742

Table 1: Comparison of different methods on the previous benchmark. For each dataset, the best results are shown in **Bold**. Reported results are averaged over three trials. Notations: $\downarrow \sim$ RMSE for regression task, $\uparrow \sim$ accuracy for classification task.

- SwitchTab consistently achieves optimal or near-optimal performance in most of the classification task
- In regression tasks, traditional methods like XGBoost or CatBoost still dominate and achieve the best results. However, SwitchTab remains highly competitive

Dataset size	45216	7043	452	200	12330	88310	518012
Feature size	17	20	236	783	12	54	54
Raw Feature (z)	✓	✓	✓	✓	✓	✓	✓
Salient Feature (s)	✓	✓	✓	✓	✓	✓	✓
Logistic Reg.	0.907	0.910	0.918	0.892	0.904	0.902	0.862
Random Forest	0.891	0.895	0.902	0.879	0.880	0.899	0.850
XGBoost	0.929	0.928	0.925	0.912	0.910	0.919	0.879
LightGBM	0.939	0.939	0.942	0.910	0.910	0.915	0.887
CatBoost	0.925	0.929	0.927	0.912	0.910	0.919	0.879
MLP	0.915	0.917	0.923	0.892	0.895	0.902	0.862
TabNet	0.766	-	0.510	-	0.653	-	0.610
TabTransformer	0.918	-	0.796	-	0.914	-	0.914
SABT	0.913	-	0.817	-	0.700	-	0.927
SwitchTab (Self-Sup.)	0.917	-	0.903	-	0.900	-	0.904
SwitchTab	0.942	-	0.923	-	0.922	-	0.931

Table 2: Comparison of different methods on classification task. For each method, we report three categories 1) raw features only, 2) salient features only, 3) plug and play using salient features. The best results are shown in **Bold**. Columns added with \star are multi-class classification tasks, reporting accuracy. The other results of binary classification tasks are evaluated with AUC.

- SwitchTab has remarkable performance lift in majority of the cases
- Salient features have immense value when integrated with original data as additional features to improve traditional methods

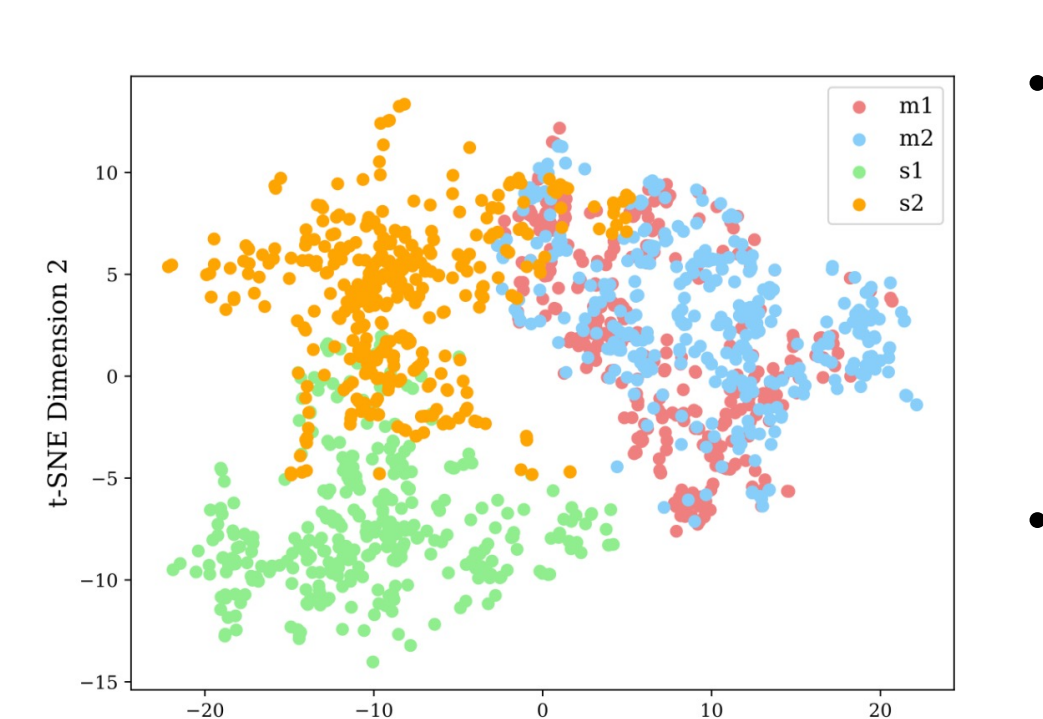


Figure 4: t-SNE visualization of mutual and salient features in two dimensional space.

- The mutual features m_1 and m_2 from SwitchTab, although heavily overlap with each other.
- The salient features s_1 and s_2 are distinctly separated, playing a dominant role in capturing the unique properties of each class and decisively contributing to the classification boundaries.