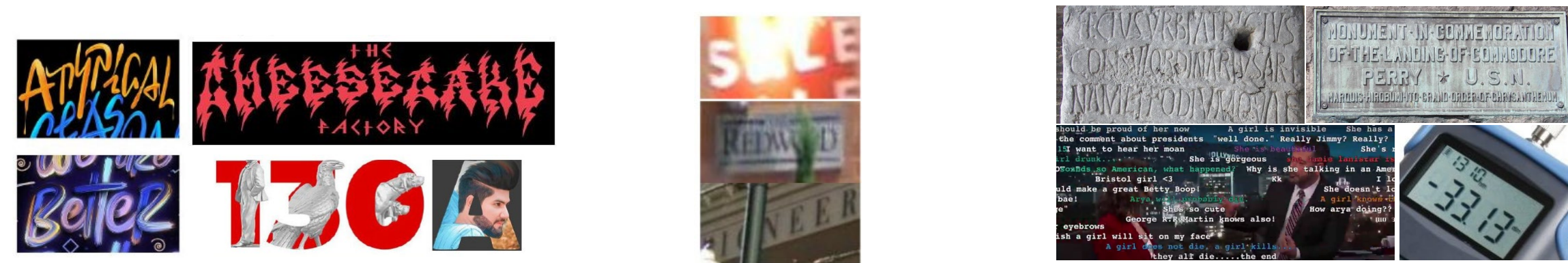




Introduction :

A. Motivation

When handling complicated text images, existing supervised text recognition methods are data-hungry.



(a) WordArt (b) Occluded text (c) More scene texts

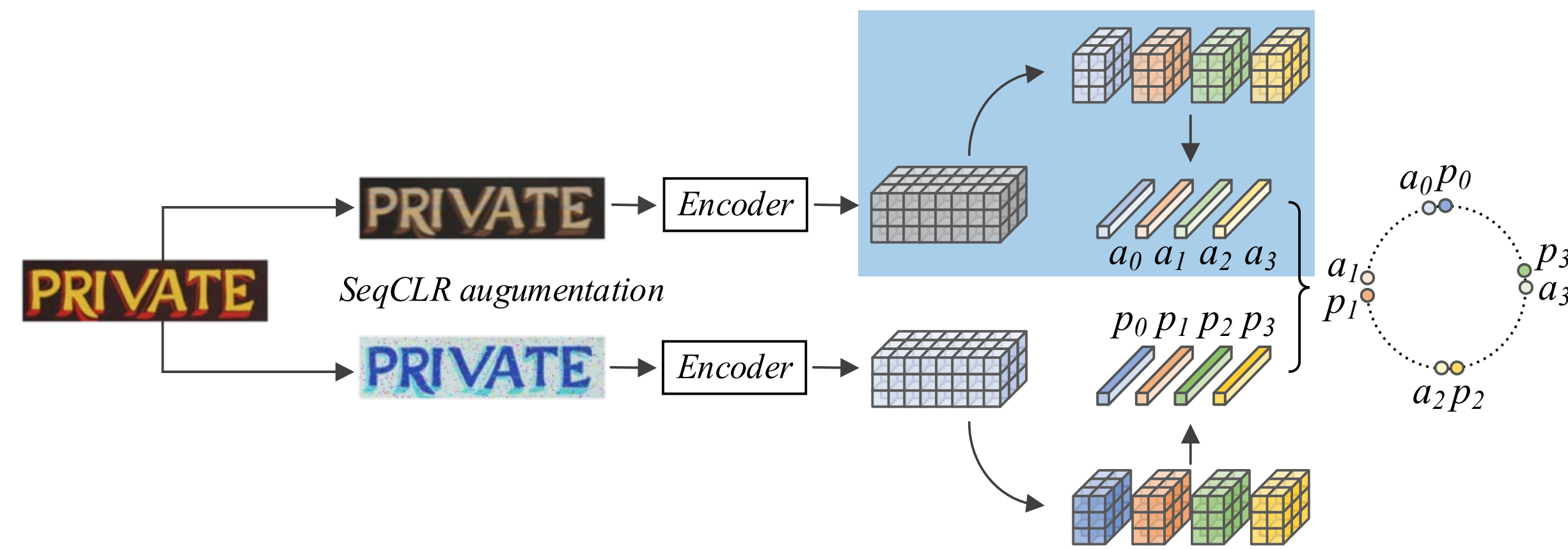
How well do recognition models generalize to arbitrary texts?

☒ supervised methods ☒ self-supervised methods

B. Objective

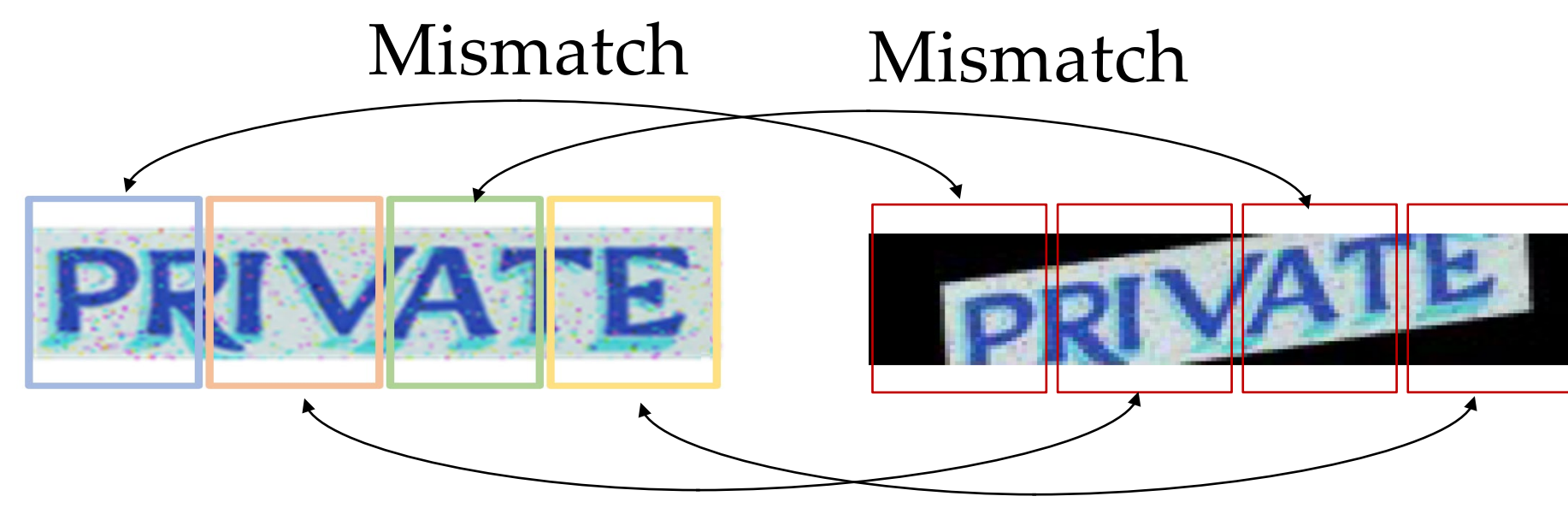
Extracting the robust visual features of characters on arbitrary texts.

Existing self-supervised pipeline



Sequence-level self-supervised learning

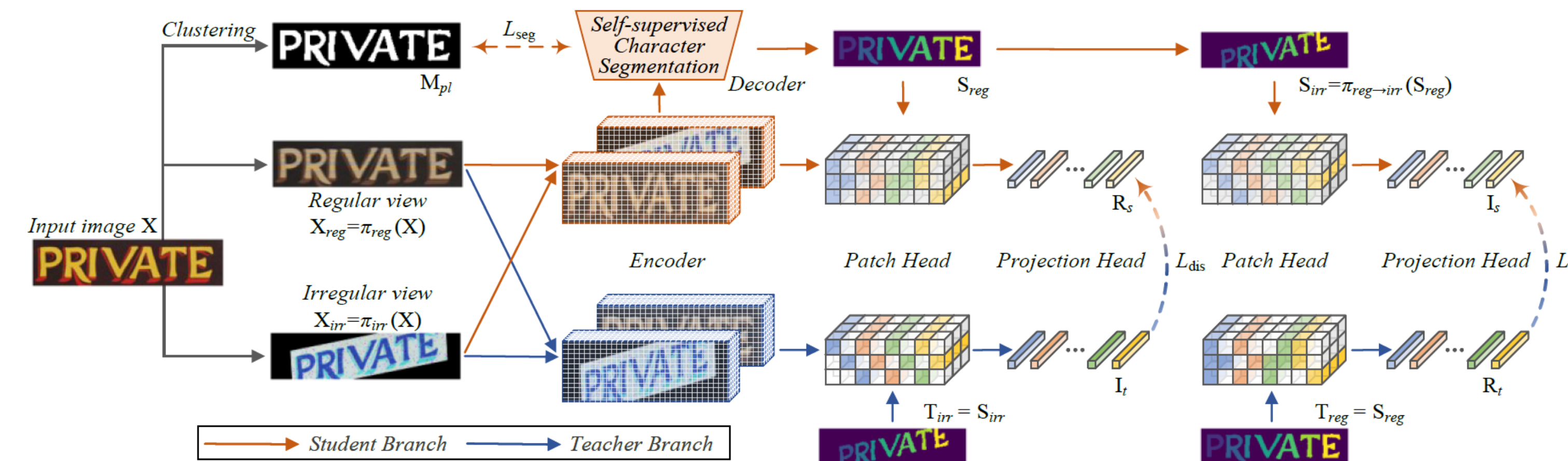
1) **Inflexible data augmentation strategy**, as large geometric transformations may cause inconsistency among the corresponding items



2) **Neglecting character structures**, which confuses networks to cause inter-character mixture

Methodology :

A. Our framework

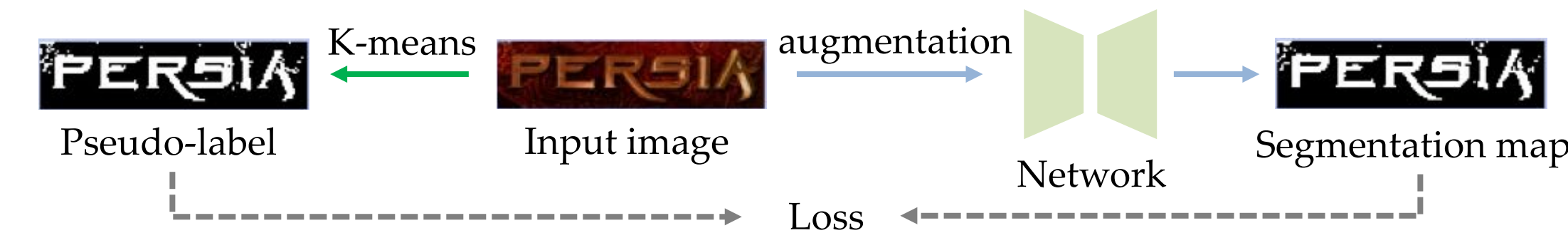


Character-level self-supervised learning

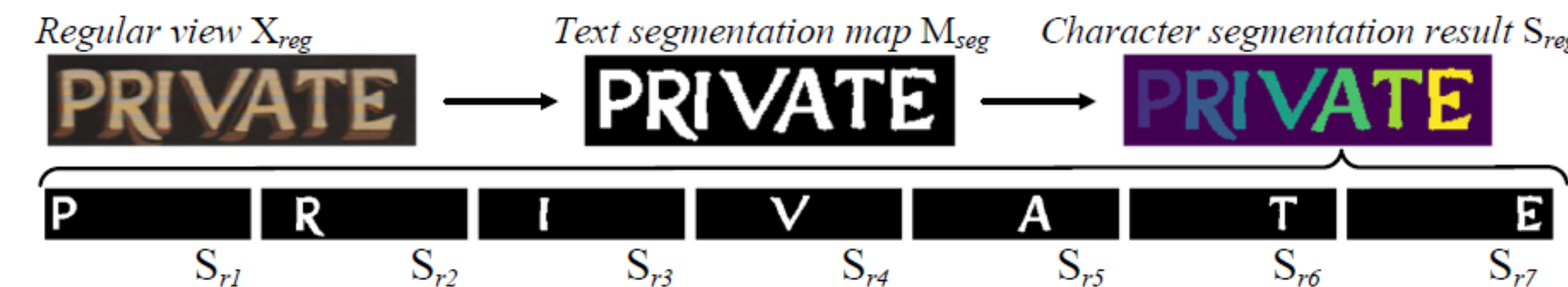
B. Details

We construct character pseudo-labels online by jointly:

(a) **self-supervised text segmentation**, producing text pseudo-labels;



(b) **clustering-based character segmentation**, producing a mask for each of its characters.



(c) **Corresponding Character Regions Alignment**

Problem:

given $\mathbf{X}_{reg} = \pi_{reg}(\mathbf{X})$, $\mathbf{X}_{irr} = \pi_{irr}(\mathbf{X})$, and \mathbf{S}_{reg} , compute $\mathbf{S}_{irr} = \pi_{reg \rightarrow irr}(\mathbf{S}_{reg})$.

Solution:

$$\pi_{reg \rightarrow irr} = \pi_{irr}(\pi_{reg}^{-1}) = \pi_{irr}$$

$$\mathbf{S}_{irr} = \pi_{irr}(\mathbf{S}_{reg})$$

(d) **Character-to-character Distillation**

- Obtain character feature representations from different views and different branches (R_s , I_s , R_t , and I_t).
- Following DINO, establish distillation loss.

Visualization :

Input image X	Regular view \mathbf{X}_{reg}	Text pseudo-label \mathbf{M}_{pl}	Character regions \mathbf{S}_{reg}
	Irregular view \mathbf{X}_{irr}	Self-supervised segmentation \mathbf{M}_{seg}	Character regions \mathbf{S}_{irr}