## Qualitative analysis

We present the T-SNE visualization of features on the STL-10 test dataset with a 40-label split in Figure 12a,12b,12c. The visualization are using trained models from FixMatch, FlexMatch, and SequenceMatch. SequenceMatch shows better feature space than FixMatch and FlexMatch with less confusing clusters.

We also visualize the T-SNE visualization of features on the SVHN test dataset and CIFAR-10 test dataset with a 40-label split in Figure 13a,13b,13c and Figure 14a,14b,14c, respectively.
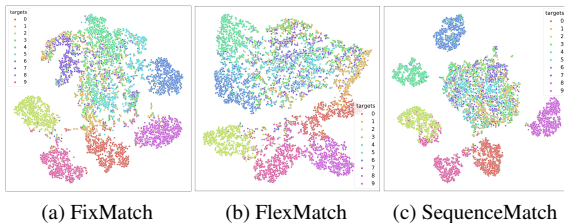


|       (a) FixMatch       |       (b) FlexMatch       |       (c) SequenceMatch       |

Figure 12. T-SNE visualization on STL-10 dataset with 40 labels.



|       (a) FixMatch       |       (b) FlexMatch       |       (c) SequenceMatch       |

Figure 13. T-SNE visualization on SVHN dataset with 40 labels.



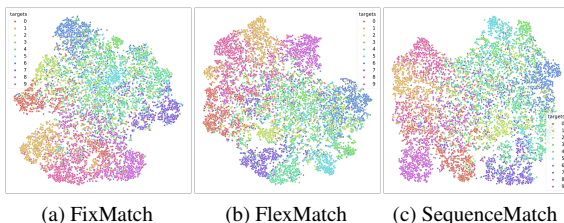|       (a) FixMatch       |       (b) FlexMatch       |       (c) SequenceMatch       |

Figure 14. T-SNE visualization on CIFAR-10 dataset with 40 labels.

## Hyperparameters setting

For reproduction, we show the detailed hyperparameter setting for each method in Table 6 and Table 7, for algorithm-dependent and algorithm-independent hyperparameters, respectively.

## Ablation study on KL loss

The additional medium augmentation requires the model to adjust to the new distribution. This is not the case with strong augmentation since the strongly augmented sample is heavily distorted, making it impossible to retrieve relevant information. As a consequence, SequenceMatch without and with KL loss obtains 5.01% and **4.80%** for CIFAR-10-40, respectively. Furthermore, after 150k iterations, the pseudo-label accuracy is 81.1%, 81.50%, and **83.20%** for FixMatch, FlexMatch, and SequenceMatch, respectively. Sequence-Match enhances the pseudo-label accuracy while improving the hard-to-learn class-wise accuracy simultaneously. This clearly demonstrates that employing medium augmentation and KL loss can reduce divergence and eliminate a confirmation bias.

## Detailed results

We also report the mean error rates of the last 20 checkpoints for various methods in Table 10. It can be seen that while most of the algorithms are overfitting to the training data at the end of the training process, our proposed method still maintains its robustness.

## ImageNet detailed results

In this section, we show the detailed results of Table 5 for ImageNet dataset on 10% labeled data. We could see that SequenceMatch outperforms previous methods in both scenarios where self-supervised pre-trained weights are included or not.

## Ablation study on medium augmentation

For the medium one, we conduct ablation studies on various types of augmentation and report the results in Table 9. We systematically test different types of augmentation for medium one, such as MoCo [15], SimCLR [8], CTAugment [3], and a combination of weak augmentation with 1 or 2 randomly selected strong augmentation. As can be seen in Table 9, using the combination of weak augmentation with 1 random strong augmentation results in the best performance.

## List of data transformations

We used the same sets of image transformations used in FixMatch [44]. For completeness, we listed all transformation operations for these augmentation strategies in Table 12.

We visualize the weak, medium, and strong augmentation examples in Figure 15 for a better understanding of the differences among the three augmentations. As we can see, the mediumly augmented examples are different from the weakly augmented ones but they are not heavily distorted like the strongly augmented ones.

Table 6. Algorithm dependent parameters. **'F-Match'** indicates FixMatch, FlexMatch, and FreeMatch.

| ALGORITHM | UDA | SEQUENCEMATCH | F-MATCH |
|---|---|---|---|
| UNLABELED DATA TO LABELED DATA RATIO (CIFAR-10/100, STL-10, SVHN) | 7 | 7 | 7 |
| UNLABELED DATA TO LABELED DATA RATIO (IMAGENET) | - | 1 | 1 |
| PRE-DEFINED THRESHOLD (CIFAR-10/100, STL-10, SVHN) | 0.8 | 0.95 | 0.95 |
| PRE-DEFINED THRESHOLD (IMAGENET) | - | 0.7 | 0.7 |
| TEMPERATURE | 0.4 | 0.5 | - |

Table 7. Algorithm independent parameters.

| DATASET | CIFAR-10 | CIFAR-100 | STL-10 | SVHN | IMAGENET |
|---|---|---|---|---|---|
| MODEL | WRN-28-2 | WRN-28-8 | WRN-37-2 | WRN-28-2 | RESNET-50 |
| WEIGHT DECAY | 5E-4 | 1E-3 | 5E-4 | 5E-4 | 3E-4 |
| BATCH SIZE | | | 64 | | 128 |
| LEARNING RATE | | | 0.03 | | |
| SGD MOMENTUM | | | 0.9 | | |
| EMA MOMENTUM | | | 0.999 | | |
| UNSUPERVISED LOSS WEIGHT | | | 1 | | |

Table 8. KL loss and results on CIFAR-10 with 40-label split.

| | WITH KL LOSS | TOP-1 |
|---|---|---|
| SEQUENCEMATCH | × | 5.01 |
| | √ | 4.80 |

Table 9. Augmentation results on CIFAR-10 with 40-label split.

| METHOD | TOP-1 |
|---|---|
| WEAK AUGMENTATION + 1 RANDOM STRONG AUGMENTATION | **4.80** |
| WEAK AUGMENTATION + 2 RANDOM STRONG AUGMENTATION | 4.91 |
| MOCO AUGMENTATION [15] | 5.97 |
| CTAUGMENT [3] | 4.85 |
| SIMCLR AUGMENTATION [8] | 5.32 |

Table 10. Mean error rates of last 20 checkpoints of all methods. There are 1000 iterations between every two checkpoints.

| DATASET | CIFAR-10 | | | CIFAR-100 | | | SVHN | | | STL-10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| # LABEL | 40 | 250 | 4000 | 400 | 2500 | 10000 | 40 | 250 | 1000 | 40 | 1000 |
| UDA | $10.65_{\pm4.97}$ | $5.67_{\pm0.28}$ | $4.58_{\pm0.07}$ | $99.0_{\pm0.0}$ | $99.0_{\pm0.0}$ | $99.0_{\pm0.0}$ | $2.5_{\pm0.54}$ | $2.06_{\pm0.02}$ | $2.01_{\pm0.03}$ | $90.0_{\pm0.0}$ | $34.88_{\pm38.98}$ |
| MPL | $8.2_{\pm1.9}$ | $8.7_{\pm1.22}$ | $4.78_{\pm0.03}$ | $48.72_{\pm0.46}$ | $29.02_{\pm0.46}$ | $22.39_{\pm0.38}$ | $11.06_{\pm6.45}$ | $2.45_{\pm0.08}$ | $2.29_{\pm0.06}$ | $44.63_{\pm7.16}$ | $7.51_{\pm0.19}$ |
| MIXMATCH | $51.5_{\pm17.51}$ | $22.14_{\pm2.83}$ | $66.57_{\pm5.38}$ | $95.87_{\pm0.24}$ | $97.88_{\pm0.36}$ | $99.0_{\pm0.0}$ | $48.86_{\pm14.71}$ | $10.16_{\pm2.7}$ | $30.09_{\pm2.42}$ | $64.99_{\pm2.3}$ | $59.5_{\pm2.48}$ |
| REMIXMATCH | $8.5_{\pm0.6}$ | $6.59_{\pm0.18}$ | $4.97_{\pm0.13}$ | $42.1_{\pm1.35}$ | $26.19_{\pm0.15}$ | $\mathbf{20.57_{\pm0.14}}$ | $21.41_{\pm12.26}$ | $10.69_{\pm0.73}$ | $11.44_{\pm1.91}$ | $34.12_{\pm5.27}$ | $6.99_{\pm0.08}$ |
| FIXMATCH | $12.85_{\pm4.51}$ | $5.26_{\pm0.08}$ | $4.43_{\pm0.02}$ | $48.87_{\pm2.48}$ | $28.84_{\pm0.4}$ | $22.93_{\pm0.14}$ | $3.5_{\pm1.05}$ | $2.06_{\pm0.01}$ | $2.11_{\pm0.02}$ | $46.71_{\pm5.25}$ | $6.14_{\pm0.25}$ |
| FLEXMATCH | $5.53_{\pm0.28}$ | $5.24_{\pm0.08}$ | $4.49_{\pm0.04}$ | $47.56_{\pm2.68}$ | $27.62_{\pm0.11}$ | $22.88_{\pm0.17}$ | $18.55_{\pm8.46}$ | $19.17_{\pm4.63}$ | $12.93_{\pm1.85}$ | $51.15_{\pm15.35}$ | $6.34_{\pm0.37}$ |
| DASH | $9.96_{\pm3.45}$ | $5.38_{\pm0.29}$ | $4.6_{\pm0.12}$ | $50.37_{\pm1.77}$ | $28.61_{\pm0.38}$ | $22.85_{\pm0.15}$ | $5.39_{\pm2.05}$ | $2.08_{\pm0.02}$ | $2.16_{\pm0.09}$ | $44.19_{\pm6.07}$ | $6.61_{\pm0.52}$ |
| COMATCH | $7.2_{\pm1.77}$ | $5.64_{\pm0.17}$ | $4.52_{\pm0.24}$ | $60.43_{\pm8.27}$ | $31.41_{\pm0.2}$ | $23.94_{\pm0.28}$ | $13.63_{\pm5.27}$ | $3.16_{\pm0.95}$ | $2.1_{\pm0.03}$ | $\mathbf{17.88_{\pm5.09}}$ | $6.07_{\pm0.02}$ |
| SIMMATCH | $5.55_{\pm0.03}$ | $5.51_{\pm0.06}$ | $4.64_{\pm0.06}$ | $42.17_{\pm0.62}$ | $30.2_{\pm0.17}$ | $23.77_{\pm0.13}$ | $14.4_{\pm0.31}$ | $3.89_{\pm2.45}$ | $2.15_{\pm0.05}$ | $27.95_{\pm6.5}$ | $6.39_{\pm0.56}$ |
| ADAMATCH | $5.33_{\pm0.22}$ | $5.34_{\pm0.05}$ | $4.71_{\pm0.02}$ | $\mathbf{40.19_{\pm1.63}}$ | $28.08_{\pm0.39}$ | $22.91_{\pm0.18}$ | $11.49_{\pm3.52}$ | $2.22_{\pm0.06}$ | $2.12_{\pm0.08}$ | $36.46_{\pm5.49}$ | $6.43_{\pm0.13}$ |
| **SEQUENCEMATCH** | $\mathbf{5.03_{\pm0.11}}$ | $\mathbf{5.07_{\pm0.11}}$ | $\mathbf{4.43_{\pm0.02}}$ | $44.52_{\pm1.01}$ | $\mathbf{27.16_{\pm0.23}}$ | $22.90_{\pm0.16}$ | $\mathbf{2.01_{\pm0.43}}$ | $\mathbf{1.89_{\pm0.01}}$ | $\mathbf{1.86_{\pm0.01}}$ | $40.21_{\pm6.11}$ | $\mathbf{5.88_{\pm0.14}}$ |

# A. Algorithm

We present the complete algorithm for SequenceMatch in Algorithm 1.

Table 11. Accuracy results on ImageNet with 10% labeled examples using [24] and [54] source code.

| Self-supervised Pre-training | Method | Top-1 | Top-5 | Params (train/test) | Epochs |
|---|---|---|---|---|---|
| None | FixMatch | 71.5 | 89.1 | 25.6M/25.6M | $\sim 300$ |
| MoCo-EMAN [6] | FixMatch-EMAN [6] | 74.0 | 90.9 | 30.0M/25.6M | $\sim 1100$ |
| None | CoMatch [24] | 73.6 | 91.6 | 30.0M/25.6M | $\sim 400$ |
| MoCo V2 [9] | CoMatch [24] | 73.7 | 91.4 | 30.0M/25.6M | $\sim 1200$ |
| None | SimMatch [54] | 74.4 | 91.6 | 30.0M/25.6M | $\sim 400$ |
| **None** | **SequenceMatch** | **75.2** | **91.9** | **25.6M/25.6M** | $\sim 400$ |

Table 12. List of transformations used in RandAugment

| Transformation | Description | Parameter | Range |
|---|---|---|---|
| Autocontrast | Maximizes the image contrast by setting the darkest (lightest) pixel to black (white). | | |
| Brightness | Adjusts the brightness of the image. $B = 0$ returns a black image, $B = 1$ returns the original image. | $B$ | [0.05, 0.95] |
| Color | Adjusts the color balance of the image like in a TV. $C = 0$ returns a black & white image, $C = 1$ returns the original image. | $C$ | [0.05, 0.95] |
| Contrast | Controls the contrast of the image. A $C = 0$ returns a gray image, $C = 1$ returns the original image. | $C$ | [0.05, 0.95] |
| Equalize | Equalizes the image histogram. | | |
| Identity | Returns the original image. | | |
| Posterize | Reduces each pixel to $B$ bits. | $B$ | [4, 8] |
| Rotate | Rotates the image by $\theta$ degrees. | $\theta$ | [-30, 30] |
| Sharpness | Adjusts the sharpness of the image, where $S = 0$ returns a blurred image, and $S = 1$ returns the original image. | $S$ | [0.05, 0.95] |
| Shear_X | Shears the image along the horizontal axis with rate $R$. | $R$ | [-0.3, 0.3] |
| Shear_Y | Shears the image along the vertical axis with rate $R$. | $R$ | [-0.3, 0.3] |
| Solarize | Inverts all pixels above a threshold value of $T$. | $T$ | [0, 1] |
| Translate_X | Translates the image horizontally by ($\lambda \times$image width) pixels. | $\lambda$ | [-0.3, 0.3] |
| Translate_Y | Translates the image vertically by ($\lambda \times$image height) pixels. | $\lambda$ | [-0.3, 0.3] |

## B. Long-tailed issue

To further prove the effectiveness of SequenceMatch, we evaluate SequenceMatch on the imbalanced SSL setting. We conduct experiments on CIFAR-10-LT, SVHN-LT, and CIFAR-100-LT with different imbalance ratios. Following [23,37,48], we use WRN-28-2 as the backbone. We consider long-tailed (LT) imbalance, where the number of data points exponentially decreases from the first class to the last, i.e., $N_k = N_1 \times \lambda^{-\frac{k-1}{L-1}}$, where $\lambda = \frac{N_1}{N_k}$. For CIFAR-10, we set $\lambda = 100, N_1 = 1000$, and $\beta = 10\%, 20\%$, and $30\%$, respectively. Similarly, we set $\lambda = 100, N_1 = 1000$, and $\beta = 20\%$ for SVHN. And for CIFAR-100, we set $\lambda =$ 20, $N_1 = 200$, and $\beta = 40$. The results are recorded in Table 13 with an average of three different runs.

Surprisingly, SequenceMatch boosts the performance by a large margin when used with ABC [23]. With an accuracy of 85.4%, SequenceMatch outperforms ABC with an 8.2% improvement when $\beta$ equals 10%.
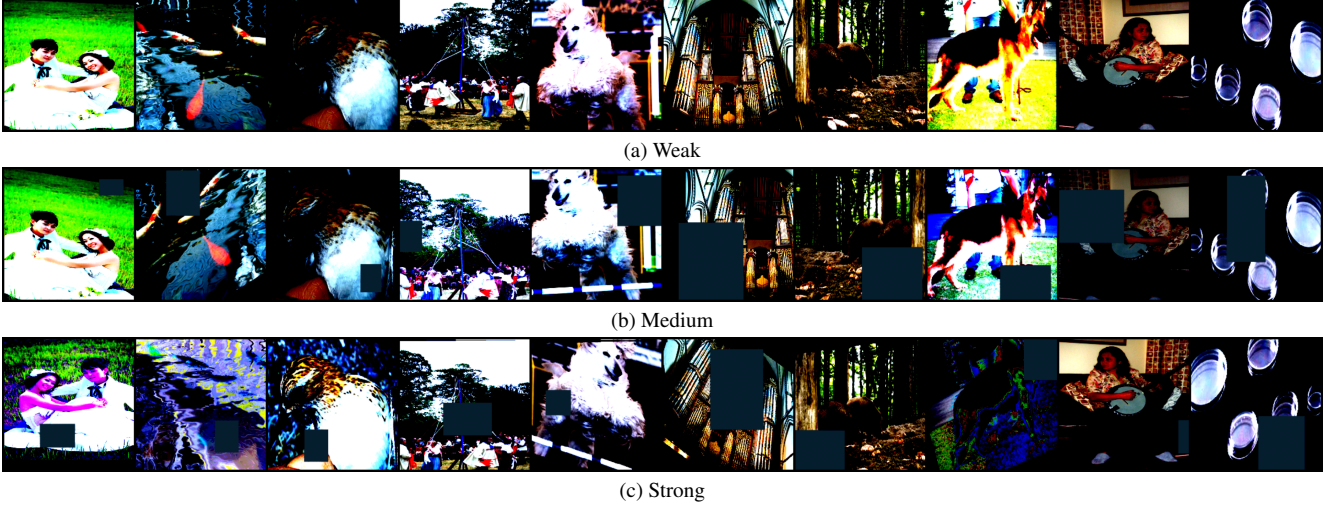
(a) Weak



(b) Medium



(c) Strong

Figure 15. Weak, medium, and strong augmented examples.

---

**Algorithm 1:** SequenceMatch algorithm

---

**Input:** Labeled batch $\mathcal{X} = (x_b, p_b) : b \in (1, \ldots, B)$, unlabeled batch $\mathcal{U} = u_b : b \in (1, \ldots, \mu B)$, confidence threshold $\tau$, unlabeled data ratio $\mu$, unlabeled loss weight $\lambda_u$, temperature $\mathbf{T}$, $\Omega$ is $\mathcal{A}_w, \mathcal{A}_m$ or $\mathcal{A}_s$

/* Cross-entropy loss for labeled data                                              */

1  $\mathcal{L}_s^{CE} = \frac{1}{B} \sum_{b=1}^{B} \mathrm{H}\left(p_b, \mathcal{A}_w\left(x_b\right)\right)$ **for** $b = 1$ **to** $\mu B$ **do**

2      $q_b\left(y \mid \Omega\left(u_b\right)\right) = p_m\left(y \mid \Omega\left(u_b\right)\right)$ // Compute prediction after applying weak data augmentation of $u_b$

3      $q_s = \frac{\exp\left(q_b/\mathbf{T}\right)}{\sum_k \exp\left(q_k/\mathbf{T}\right)}$. // Sharpen the output probability

/* Cross-entropy loss with pseudo-label and confidence threshold for high-confidence unlabeled   */

4  $\mathcal{L}_u^{\mathrm{CE}} = \frac{1}{\mu B} \sum_{b=1}^{\mu B}\left(\mathbb{1}\left(\max\left(q_b^w\right) \geq \tau\right) \mathrm{H}\left(\hat{q}_b, p_m\left(y \mid \mathcal{A}_s\left(u_b\right)\right)\right) + \mathbb{1}\left(\max\left(q_b^w\right) < \tau\right) \mathrm{H}\left(q_s^w \mid q_b\left(\mathcal{A}_s\left(u_b\right)\right)\right)\right)$

/* Kullback-Leibler divergence loss with each pair of augmented examples                 */

5  $\mathcal{L}_{\mathrm{KL}}^{w-m} = \frac{1}{\mu B} \sum_{b=1}^{\mu B} \mathbb{1}\left(\max\left(q_b^w\right) \geq \tau\right) \qquad D_{\mathrm{KL}}\left(q_s^w \mid p_m\left(y \mid \mathcal{A}_m\left(u_b\right)\right)\right)$

6  $\mathcal{L}_{\mathrm{KL}}^{m-s} = \frac{1}{\mu B} \sum_{b=1}^{\mu B} \mathbb{1}\left(\max\left(q_b^m\right) \geq \tau\right) \qquad D_{\mathrm{KL}}\left(q_s^m \mid p_m\left(y \mid \mathcal{A}_s\left(u_b\right)\right)\right)$

7  $\mathcal{L}_{\mathrm{KL}}^{w-s} = \frac{1}{\mu B} \sum_{b=1}^{\mu B} \mathbb{1}\left(\max\left(q_b^w\right) \geq \tau\right) \qquad D_{\mathrm{KL}}\left(q_s^w \mid p_m\left(y \mid \mathcal{A}_s\left(u_b\right)\right)\right)$

8  $\mathcal{L}_u = \mathcal{L}_u^{\mathrm{CE}} + \mathcal{L}_{\mathrm{KL}}^{w-m} + \mathcal{L}_{\mathrm{KL}}^{m-s} + \mathcal{L}_{\mathrm{KL}}^{w-s}$

9  **return** $\mathcal{L}_s^{CE} + \lambda_u \mathcal{L}_u$

---

Table 13. Overall accuracy under the long-tailed setting

| | CIFAR-10-LT | | | SVHN-LT | CIFAR-100-LT |
|---|---|---|---|---|---|
| | | $\lambda = 100$ | | $\lambda = 100$ | $\lambda = 20$ |
| ALGORITHM | $\beta = 10\%$ | $\beta = 20\%$ | $\beta = 30\%$ | $\beta = 20\%$ | $\beta = 40\%$ |
| VANILLA | - | $55.3_{\pm 1.30}$ | - | $77.0_{\pm 0.67}$ | $40.1_{\pm 1.15}$ |
| VAT [30] | - | $55.3_{\pm 0.88}$ | - | $81.3_{\pm 0.47}$ | $40.4_{\pm 0.34}$ |
| BALMS [37] | - | $70.7_{\pm 0.59}$ | - | $87.6_{\pm 0.53}$ | $50.2_{\pm 0.54}$ |
| FIXMATCH [44] | $70.0_{\pm 0.59}$ | $72.3_{\pm 0.33}$ | $74.9_{\pm 0.63}$ | $88.0_{\pm 0.30}$ | $51.0_{\pm 0.20}$ |
| W/ CRES T+PDA [48] | $73.9_{\pm 0.40}$ | $76.6_{\pm 0.46}$ | $74.9_{\pm 0.63}$ | $89.1_{\pm 0.69}$ | $51.6_{\pm 0.29}$ |
| W/ DARP [19] | - | $73.7_{\pm 0.98}$ | - | $88.6_{\pm 0.19}$ | $51.4_{\pm 0.37}$ |
| W/ DARP+CRT [19] | $74.6_{\pm 0.98}$ | $78.1_{\pm 0.89}$ | $77.6_{\pm 0.73}$ | $89.9_{\pm 0.44}$ | $54.7_{\pm 0.46}$ |
| W/ ABC [23] | $77.2_{\pm 1.60}$ | $81.1_{\pm 0.82}$ | $81.5_{\pm 0.29}$ | $92.0_{\pm 0.38}$ | $56.3_{\pm 0.19}$ |
| **SEQUENCEMATCH** | $\mathbf{85.4_{\pm 0.01}}$ | $\mathbf{81.5_{\pm 0.75}}$ | $\mathbf{82.2_{\pm 0.25}}$ | $\mathbf{92.4_{\pm 0.06}}$ | $\mathbf{57.2_{\pm 0.09}}$ |