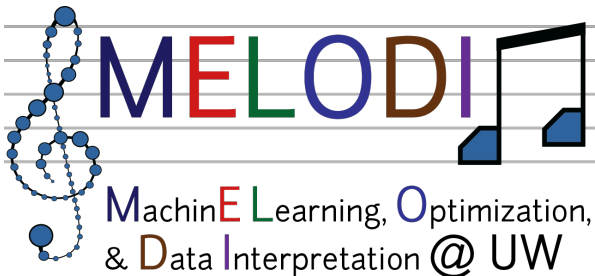


# Time-Consistent Self-Supervision for Semi-Supervised Learning

Tianyi Zhou\*, Shengjie Wang\*, Jeff A. Bilmes

*University of Washington, Seattle*



**Can SSL achieve fully-supervision's accuracy using similar amount of computation?**

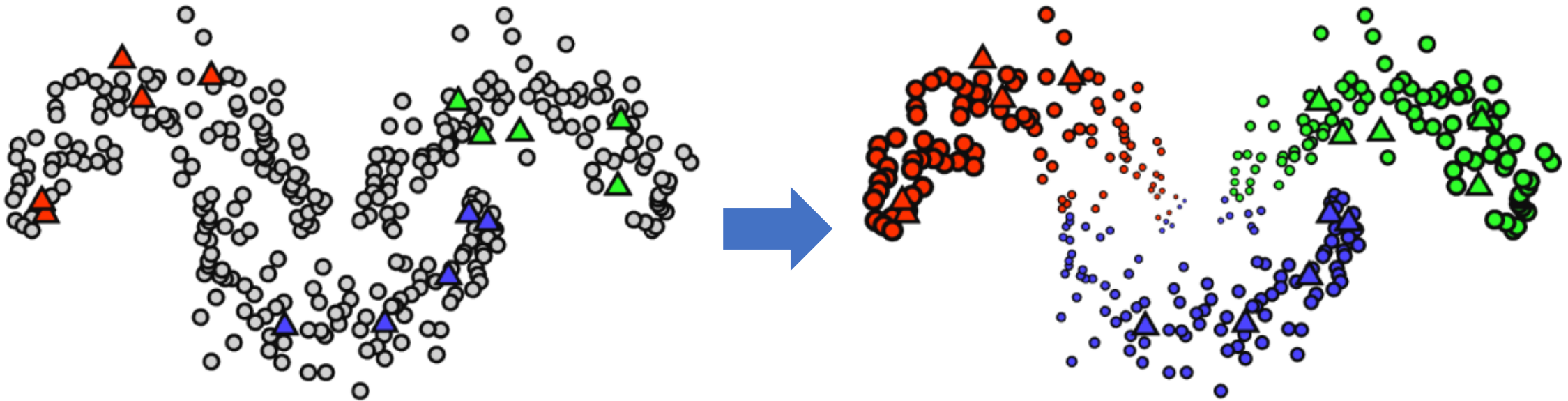
**Yes!**

**How?**

**Select unlabeled data with time-consistent prediction for self-supervision.**

# Semi-Supervised Learning with **Spatial Consistency**

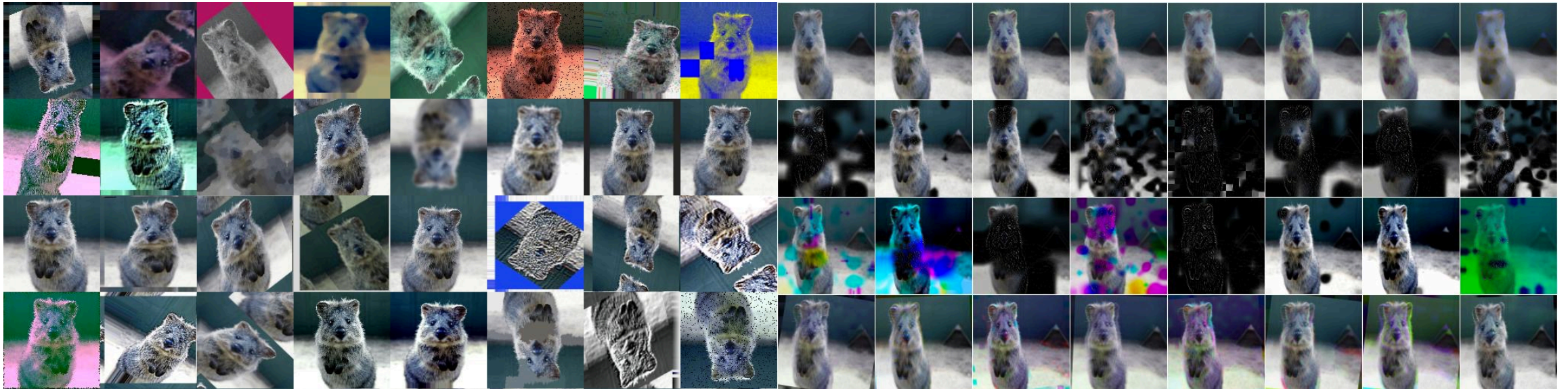
- Idea: Samples with similar features/embeddings have similar labels
- Previous: Label/measure propagation; manifold regularization
- More recent: The same idea inspires graphical neural networks



credit: [Isken et al. 2019]

# Semi-Supervised Learning with **Pseudo Targets**

- Idea: average the model output of an unlabeled sample over multiple augmentations/steps; use the average as training target.
- Fit in Deep learning: encourage spatial consistency around single sample; working with data augmentation and inductive bias of DNNs.
- Drawbacks: can be wrong on some samples; early-stage model is poor
  - **In practice**: select samples with high confidence, but DNNs can be over-confident.



# A Recipe of **Self-Supervision** on Unlabeled Data

- Consistency loss:

An unlabeled sample and its augmentation should have **similar predictions**.

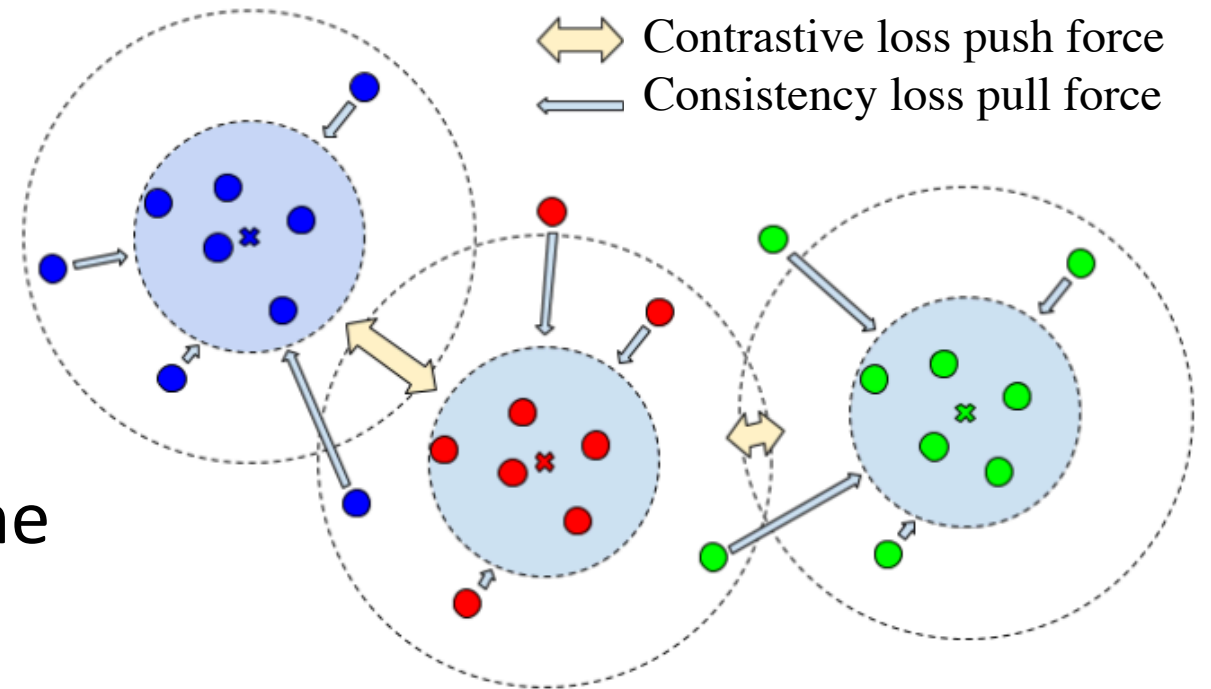
- Contrastive loss/Triplet loss:

Different samples (and their augmentations) should have **more different predictions** than the same sample and its augmentations.

- Cross Entropy loss defined on **pseudo targets**.

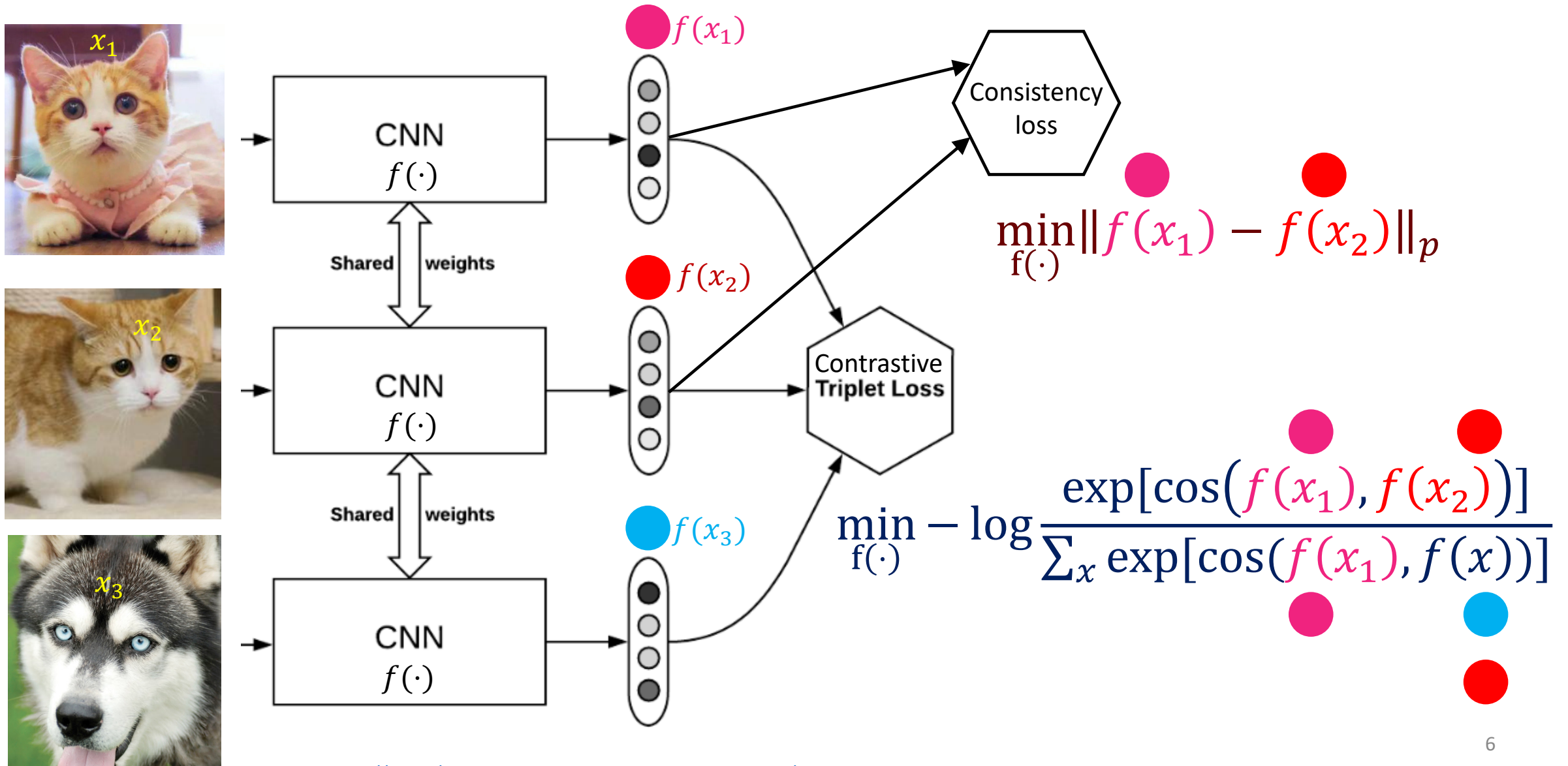
- Our SSL objective combines the three losses.

*Each color represent a sample and its augmentations*



credit: [Brabandere et al. 2017]

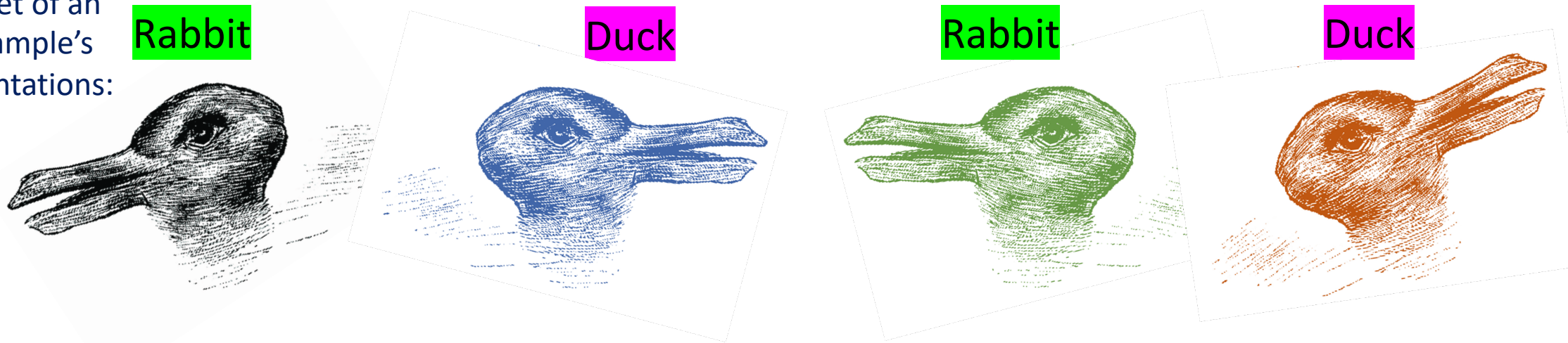
# An Example of Consistency/Contrastive Loss



# Problem of Current SSL: **time-inconsistency**

- The pseudo target depends on model in-training and is time-variant.
- Hence, the training objective is time-inconsistent!
- DNN is confusing itself in self-supervision.
- Possible outcomes: divergence, concept drift, catastrophic forgetting, etc.

Pseudo Target of an unlabeled sample's data augmentations:



Training Time

**Self-supervision losses depend on pseudo targets (or model outputs), which should be time-consistent!**



# Time-Consistency (TC)

- We select unlabeled data with consistent predictions/outputs for self-supervision in SSL by using a curriculum.
- (instantaneous) Time consistency of sample  $x$  at step- $t$  (e.g.,  $t^{\text{th}}$  mini-batch):

$$a^t(x) \triangleq D_{KL}(p^{t-1}(x) || p^t(x)) + \left| \log \frac{p^{t-1}(x)[y^{t-1}(x)]}{p^t(x)[y^{t-1}(x)]} \right|$$

$p^t(x)$ : output distribution over classes for  $x$  at step- $t$

$y^t(x)$ : predicted class for  $x$  at step- $t$

# Time-Consistency (TC)

- (instantaneous) Time consistency of  $x$  at step- $t$ :

$$a^t(x) \triangleq D_{KL}(p^{t-1}(x) || p^t(x)) + \left| \log \frac{p^{t-1}(x)[y^{t-1}(x)]}{p^t(x)[y^{t-1}(x)]} \right|$$

○  $y^t(x) = \arg \max_i p^t(x)[i]$ , i. e., the class with the highest probability.

- 1<sup>st</sup> term: *KL-divergence* between the predictions at step  $t$  and  $t-1$ .
- 2<sup>nd</sup> term: change of confidence on the predicted class between step  $t$  and  $t-1$ .

# Time-Consistency (TC)

- (instantaneous) Time consistency of  $x$  at step- $t$ :

$$a^t(x) \triangleq D_{KL}(p^{t-1}(x) || p^t(x)) + \left| \log \frac{p^{t-1}(x)[y^{t-1}(x)]}{p^t(x)[y^{t-1}(x)]} \right|$$

○  $y^t(x) = \arg \max_i p^t(x)[i]$ , i. e., the class with the highest probability.

- 1<sup>st</sup> term: *KL-divergence* between the predictions at step  $t$  and  $t-1$ .
- 2<sup>nd</sup> term: change of confidence on the predicted class between step  $t$  and  $t-1$ .

# Time-Consistency (TC)

- (instantaneous) Time consistency of  $x$  at step- $t$ :

$$a^t(x) \triangleq D_{KL}(p^{t-1}(x) || p^t(x)) + \left| \log \frac{p^{t-1}(x)[y^{t-1}(x)]}{p^t(x)[y^{t-1}(x)]} \right|$$

- Time Consistency (TC): smooth  $-a^t(x)$  by exponential moving average over time steps:

$$c^t(x) = \gamma_c(-a^t(x)) + (1 - \gamma_c)c^{t-1}(x)$$

Time-Consistency relates to

## Catastrophic Forgetting in Training Dynamics

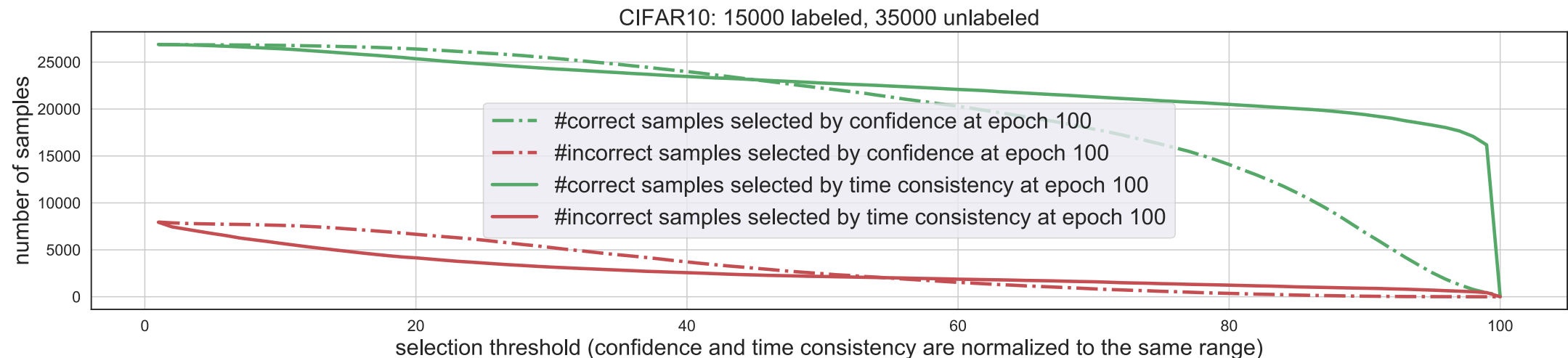
- $a^t(x')$  is an upper-bound on the forgetfulness of catastrophic forgetting on labeled data if adding an unlabeled sample  $x'$  and its pseudo targets to training:

$$\text{Forgetfulness} \triangleq \frac{1}{\eta} \left| \sum_{x \in \mathcal{L}} \left[ \ell(x; \theta^{t+1}) - \ell(x; \hat{\theta}^{t+1}) \right] \right|$$

- $\ell(x; \theta)$ : loss of model  $\theta$  on sample  $x$ ;
  - Assume the loss on labeled data  $L$  is close to 0 after warm-starting epochs, i.e.,  $\sum_{x \in L} \ell(x; \theta^t) \approx 0$ .
  - $\theta^t$ : model-at-step- $t$  updated by **labeled data**;
  - $\hat{\theta}^t$ : model-at-step- $t$  updated by **labeled data +  $x'$** ;
- A small  $a^t(x')$  means adding  $x'$  and its pseudo target to training does not cause forgetting of labeled data (and previously trained unlabeled-data).

# Empirical Evidence of Time Consistency

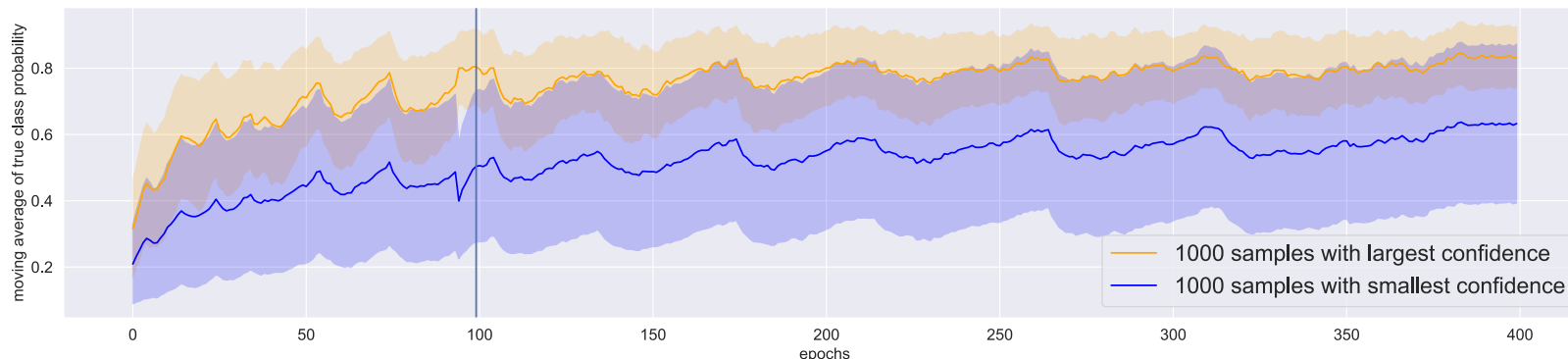
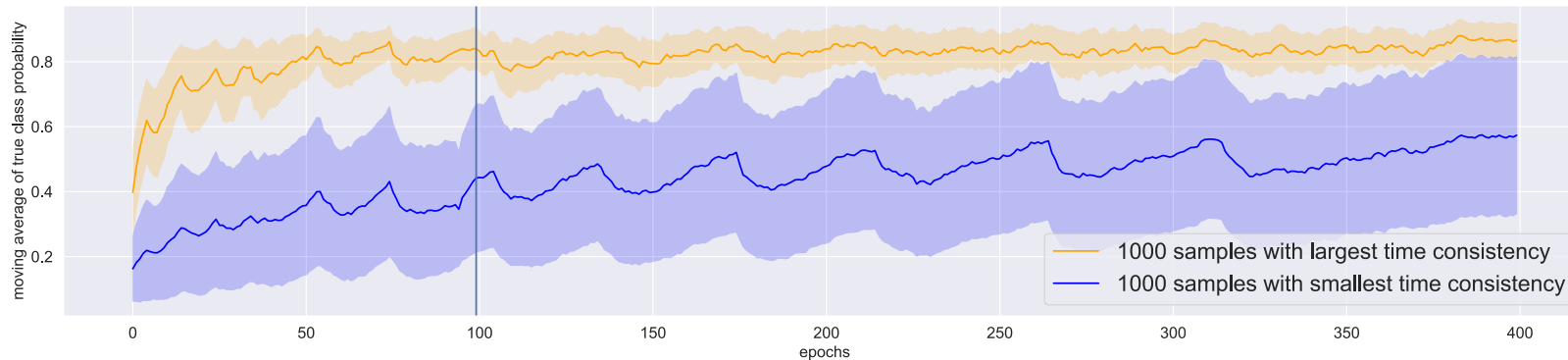
- Split CIFAR10 training set into two subsets of 15000 and 35000 samples.
- Train WideResNet-18-2 on the 15000 samples, test it on the 35000 samples.
- Time consistency performs better than confidence in identifying the unlabeled samples correctly predicted by the current model.



Computed time-consistency and confidence at epoch 100 of training WideResNet-28-2. The x-axis shows the validation samples selected using different thresholds on the two metrics (normalized to  $[0, 100]$ ). The y-axis reports correct v.s. incorrect predictions over the selected samples.

# Persistence of Time Consistency

- Time consistency performs better in predicting the future dynamics, i.e., it identifies samples whose predictions stay correct stably in the future.



- Computed time-consistency (top) and confidence (bottom) at epoch 100 of training WideResNet-28-2 on CIFAR10.
- Select the top 1000 and bottom 1000 validation samples based on the two metrics.
- Compare the moving average of true class probability of the selected samples across epochs.

# TC-SSL Algorithm

- In each step, select unlabeled samples with large time-consistency and optimize our SSL objective on them.
- Add warm-start epochs and apply exponential weighted sampling to encourage exploration in early stages.
- Remove samples with extremely high confidence since they contribute nearly zero gradients.
- Follow previous works: Mix-Up, sharpen predicted probability as pseudo target, duplicate labeled data to similar amount of selected unlabeled data, etc.

---

## Algorithm 1 Time-Consistent SSL (TC-SSL)

---

```

1: input:  $\mathcal{U}, \mathcal{L}, \pi(\cdot; \eta), \eta^{1:T}, f(\cdot; \theta), G(\cdot)$ ;
2: hyperparameters:  $T_0, T, \lambda_{cs}, \lambda_{ct}, \lambda_{ce}, \gamma_\theta, \gamma_c, \gamma_k$ ;
3: initialize:  $\theta^0, k^1$ ;
4: for  $t \in \{1, \dots, T\}$  do
5:   if  $t \leq T_0$  then
6:      $\theta^t \leftarrow \theta^{t-1} + \pi \left( \sum_{(x,y) \in \mathcal{L}} \nabla_\theta \ell_{ce}(x, y; \theta^{t-1}); \eta^t \right)$ 
7:   else
8:      $S^t = \operatorname{argmax}_{S: S \subseteq \mathcal{U}, |S|=k^t} \sum_{x \in S} c^t(x)$  or
9:     Draw  $k^t$  samples from  $\Pr(x \in S^t) \propto \exp(c^t(x))$ ;
10:     $\theta^t \leftarrow \theta^{t-1} + \pi \left( \nabla_\theta L^t(\theta^{t-1}); \eta^t \right)$  (ref. Eq. (11));
11:   end if
12:    $p^t(x) \leftarrow \frac{\exp(f(x; \theta^t)[y])}{\sum_{y'=1}^C \exp(f(x; \theta^t)[y'])}, \forall y \in [C], x \in \mathcal{U}$ ;
13:   if  $t = 1$  then
14:      $\overline{\theta^t} \leftarrow \theta^t, c^t(x) \leftarrow 0, \forall x \in \mathcal{U}$ 
15:   else
16:     Compute  $a^t(x)$  (ref. Eq (1)),  $\forall x \in \mathcal{U}$ ;
17:   end if
18:    $c^{t+1}(x) \leftarrow \gamma_c(-a^t(x)) + (1 - \gamma_c)c^{t-1}(x), \forall x \in \mathcal{U}$ ;
19:    $\overline{\theta^{t+1}} \leftarrow \gamma_\theta \theta^t + (1 - \gamma_\theta)\overline{\theta^t}$ ;
20:    $k^{t+1} \leftarrow (1 + \gamma_k) \times k^t$ ;
21: end for

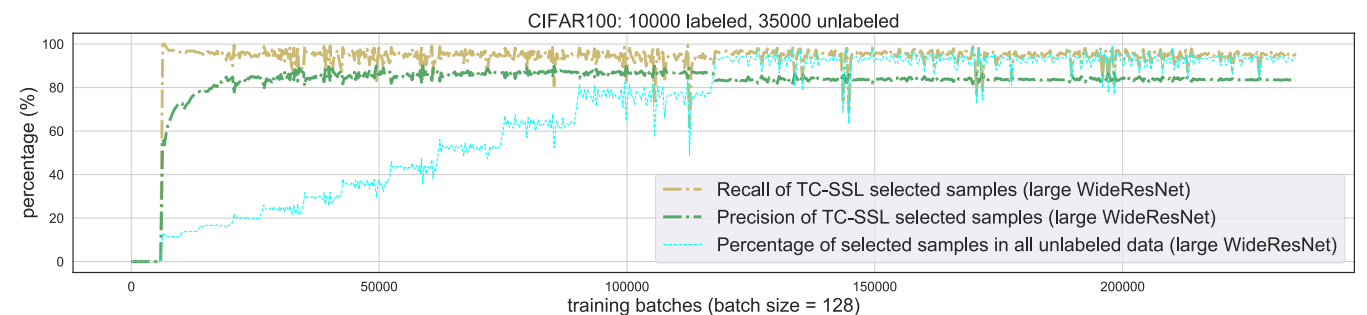
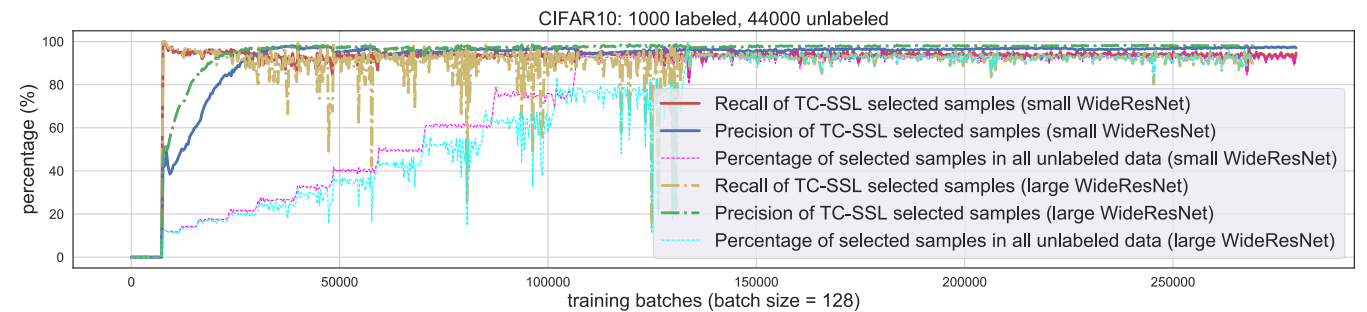
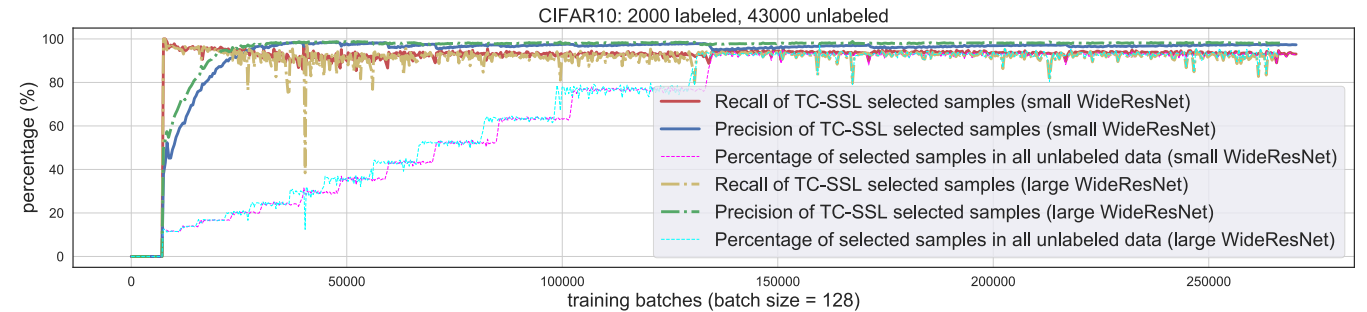
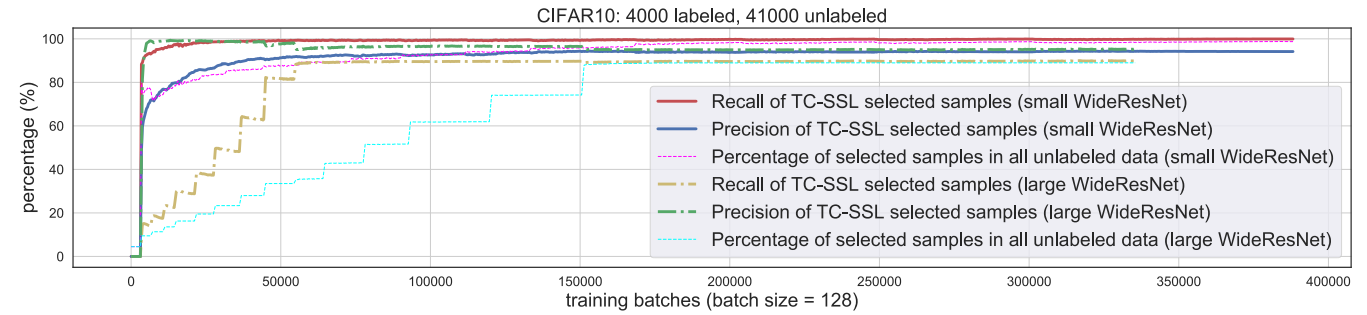
```

---



# Quality of Selected Pseudo Targets in TC-SSL

- TC-SSL produces a curriculum of unlabeled data whose pseudo targets are of high precision and recall throughout the course of training;
- TC-SSL gradually increase the use of unlabeled data rather than adding all of them to training at the very beginning.



# Experimental Results

- TC-SSL achieves SOTA performance on CIFAR10, CIFAR100, STL10 of different labeled/unlabeled splittings (more results in paper).

Table 1. Test error rate (mean $\pm$ variance) of SSL methods training a small WideResNet and a large WideResNet on **CIFAR10**. Baselines: Pseudo Label (Lee, 2013),  $\Pi$ -model (Sajjadi et al., 2016), VAT (Miyato et al., 2019), Mean Teacher (Tarvainen & Valpola, 2017), MixMatch (Berthelot et al., 2019), ReMixMatch (Berthelot et al., 2020).

Benchmark	CIFAR10 (small WideResNet-28-2)				CIFAR10 (large WideResNet-28-135)				
	labeled/unlabeled	500/44500	1000/44000	2000/43000	4000/41000	500/44500	1000/44000	2000/43000	4000/41000
Pseudo Label		40.55 $\pm$ 1.70	30.91 $\pm$ 1.73	21.96 $\pm$ 0.42	16.21 $\pm$ 0.11	-	-	-	-
$\Pi$ -model		41.82 $\pm$ 1.52	31.53 $\pm$ 0.98	23.07 $\pm$ 0.66	5.70 $\pm$ 0.13	-	-	-	-
VAT		26.11 $\pm$ 1.52	18.68 $\pm$ 0.40	14.40 $\pm$ 0.15	11.05 $\pm$ 0.31	-	-	-	-
Mean Teacher		42.01 $\pm$ 5.86	17.32 $\pm$ 4.00	12.17 $\pm$ 0.22	10.36 $\pm$ 0.25	-	-	-	-
MixMatch		9.65 $\pm$ 0.94	7.75 $\pm$ 0.32	7.03 $\pm$ 0.15	6.24 $\pm$ 0.06	8.44 $\pm$ 1.04	7.38 $\pm$ 0.63	6.51 $\pm$ 0.48	5.12 $\pm$ 0.31
ReMixMatch		-	5.73 $\pm$ 0.16	-	5.14 $\pm$ 0.04	-	-	-	-
TC-SSL (ours)		9.14 $\pm$ 0.88	6.15 $\pm$ 0.23	5.85 $\pm$ 0.10	5.07 $\pm$ 0.05	6.04 $\pm$ 0.39	3.81 $\pm$ 0.19	3.79 $\pm$ 0.21	3.54 $\pm$ 0.06

# Experimental Results

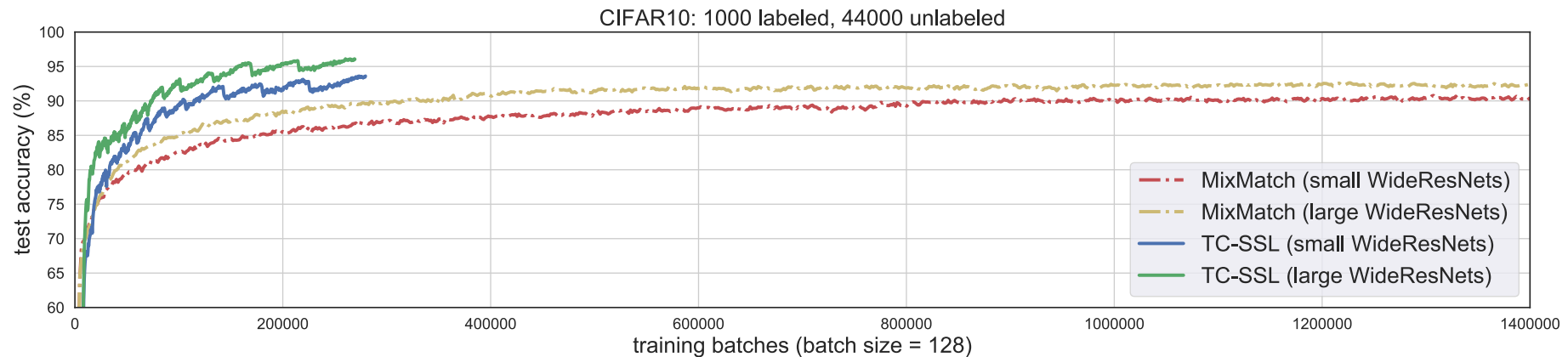
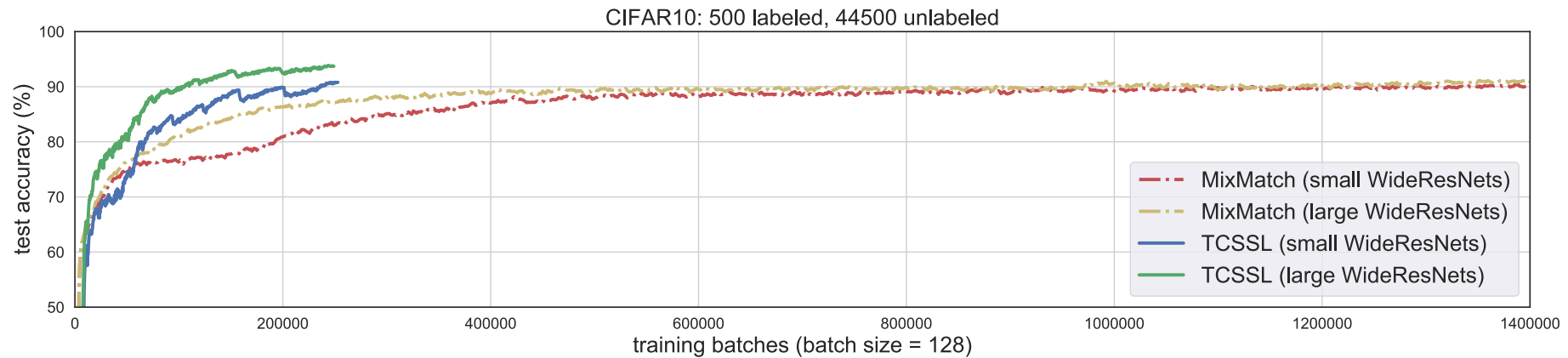
- TC-SSL achieves SOTA performance on CIFAR10, CIFAR100, STL10 of different labeled/unlabeled splittings (more results in paper).

Benchmark	STL10
labeled/unlabeled	5k/100k
CutOut (11.0M)	12.74
IIC (N/A)	11.20
MixMatch (23.8M)	5.59
TC-SSL (ours, 5.9M)	4.82

Benchmark	CIFAR100 (WideResNet-28-135)		
labeled/unlabeled	2500/42500	5000/40000	10000/35000
SWA	-	-	28.80
MixMatch	44.20 $\pm$ 1.18	34.62 $\pm$ 0.63	25.88 $\pm$ 0.30
TC-SSL (ours)	31.95 $\pm$ 0.55	26.98 $\pm$ 0.51	22.10 $\pm$ 0.37

# Experimental Results

- TC-SSL significantly improves SSL efficiency.
- It achieves high accuracy using much fewer but more informative and time-consistent training batches with more accurate pseudo targets.



# Ablation Study

- Test error rate (mean $\pm$ variance) of TC-SSL variants training WideResNet on CIFAR10;
- no consistency: TC-SSL without consistency loss;
- no contrastive: TC-SSL without contrastive loss;
- no PseudoLabel: TC-SSL without cross entropy loss for unlabeled data;
- no TC-selection: replace TC-based selection/sampling with uniform sampling.

labeled/unlabeled	500/44500	1000/44000	2000/43000	4000/41000
TC-SSL (ours)	6.04 $\pm$ 0.39	3.81 $\pm$ 0.19	3.79 $\pm$ 0.21	3.54 $\pm$ 0.06
TC-SSL (no consistency)	7.51 $\pm$ 0.56	5.31 $\pm$ 0.23	3.82 $\pm$ 0.20	3.58 $\pm$ 0.06
TC-SSL (no contrastive)	7.56 $\pm$ 0.52	5.35 $\pm$ 0.28	3.96 $\pm$ 0.25	3.66 $\pm$ 0.08
TC-SSL (no PseudoLabel)	41.05 $\pm$ 2.32	23.64 $\pm$ 1.17	14.37 $\pm$ 0.69	9.87 $\pm$ 0.22
TC-SSL (no TC-selection)	12.25 $\pm$ 0.81	6.39 $\pm$ 0.44	4.68 $\pm$ 0.35	4.05 $\pm$ 0.13

# Take-home Messages

---



Time-consistency is critical to semi-supervised learning;



We derive a novel time-consistency metric with theoretical support on avoiding catastrophic forgetting and plenty of empirical evidences;



We provide a recipe of self-supervision losses: consistency + contrastive;



TC-SSL, the proposed algorithm, achieves SOTA performance on several SSL benchmarks and considerably improves efficiency.



Thank you!

---

- For questions and discussions, please join our Q&A session.
  - July 15 Web 10:00 AM PDT
  - July 15 Web 23:00 PM PDT