

## Don't miss the Mismatch: Investigating the Objective Function Mismatch for Unsupervised Representation Learning (Supplementary Material)

### A Proofs

We measure our metrics on the mean losses during cross-validation instead of calculating the metrics for each round and taking the average. We proof that both variants are equivalent for M3 while measuring SM3 on the mean losses leads to a lower bound, given that all models converge at step  $s_n$ .

**Proposition 1** *The MM3 of the average metric value tuples  $\frac{1}{h} \sum_{0 < c \leq h} M_c^P$  and  $\frac{1}{h} \sum_{0 < c \leq h} M_c^T$  with  $0 < c \leq h$  is equivalent to the average MM3 of the individual tuples  $\frac{1}{h} \sum_{0 < c \leq h} \text{MM3}(M_c^T, M_c^P)$ , given that the tuples are measured for the same steps  $S$  and converge at the same step  $s_n$ .*

*Proof.*

$$\begin{aligned} \frac{1}{h} \sum_{0 < c \leq h} \text{MM3}(M_c^T, M_c^P) &= \frac{1}{h} \sum_{0 < c \leq h} \frac{1}{n} \sum_{0 < i \leq n} (m_{ic}^T - m_{ic}^P) \\ &= \frac{1}{n} \sum_{0 < i \leq n} \frac{1}{h} \sum_{0 < c \leq h} (m_{ic}^T - m_{ic}^P) \\ &= \frac{1}{n} \sum_{0 < i \leq n} \left( \frac{1}{h} \sum_{0 < c \leq h} m_{ic}^T - \frac{1}{h} \sum_{0 < c \leq h} m_{ic}^P \right) \\ &= \text{MM3} \left( \frac{1}{h} \sum_{0 < c \leq h} M_c^T, \frac{1}{h} \sum_{0 < c \leq h} M_c^P \right) \end{aligned}$$

**Corollary 1**  $\frac{1}{h} \sum_{0 < c \leq h} \text{M3}(m_c^T, m_c^P) = \text{M3}(\frac{1}{h} \sum_{0 < c \leq h} m_c^T, \frac{1}{h} \sum_{0 < c \leq h} m_c^P)$

**Proposition 2** *The MSM3 of the average metric value tuples  $\frac{1}{h} \sum_{0 < c \leq h} M_c^P$  and  $\frac{1}{h} \sum_{0 < c \leq h} M_c^T$  with  $0 < c \leq h$  is a lower bound of the average MSM3 of the individual tuples  $\frac{1}{h} \sum_{0 < c \leq h} \text{MSM3}(M_c^T)$ , given that the tuples are measured for the same steps  $S$  and converge at the same step  $s_n$ .*

*Proof.*

$$\begin{aligned} \frac{1}{h} \sum_{0 < c \leq h} \text{MSM3}(M_c^T) &= \frac{1}{h} \sum_{0 < c \leq h} \frac{1}{n} \sum_{0 < i \leq n} \left( m_{ic}^T - \min_{0 < j \leq i} (m_{jc}^T) \right) \\ &= \frac{1}{n} \sum_{0 < i \leq n} \frac{1}{h} \sum_{0 < c \leq h} \left( m_{ic}^T - \min_{0 < j \leq i} (m_{jc}^T) \right) \\ &= \frac{1}{n} \sum_{0 < i \leq n} \left( \frac{1}{h} \sum_{0 < c \leq h} m_{ic}^T - \frac{1}{h} \sum_{0 < c \leq h} \min_{0 < j \leq i} (m_{jc}^T) \right) \\ &\geq \frac{1}{n} \sum_{0 < i \leq n} \left( \frac{1}{h} \sum_{0 < c \leq h} m_{ic}^T - \min_{0 < j \leq i} \left( \frac{1}{h} \sum_{0 < c \leq h} m_{jc}^T \right) \right) && \text{since } \frac{1}{h} \sum_{0 < c \leq h} \min_{0 < j \leq i} (m_{jc}^T) \leq \min_{0 < j \leq i} \left( \frac{1}{h} \sum_{0 < c \leq h} m_{jc}^T \right) \\ &\geq \text{MSM3} \left( \frac{1}{h} \sum_{0 < c \leq h} M_c^T \right) \end{aligned}$$

**Corollary 2**  $\frac{1}{h} \sum_{0 < c \leq h} \text{SM3}(m_c^T, m_c^P) \geq \text{SM3}(\frac{1}{h} \sum_{0 < c \leq h} m_c^T, \frac{1}{h} \sum_{0 < c \leq h} m_c^P)$

## B Additional Model and Training Details

**CNN Encoders:** We consider a family of convolutional encoders with four Conv-BatchNorm-ReLU layers. Filter widths are  $[32, 64, 128, f]$  and paddings are "valid". For input sizes of  $32 \times 32$  (Cifar10, Cifar100), kernel sizes are  $[3, 3, 3, 2]$  and strides are  $[2, 2, 2, 1]$ ; for input sizes of  $64 \times 64$  (3dshapes, PCam), kernel sizes are  $[4, 4, 4, 3]$  and strides are  $[2, 2, 2, 2]$ . Weights are initialized with the standard TensorFlow initialization (kernel\_initializer="glorot\_uniform", bias\_initializer="zeros"). We vary  $f$  in  $[4, 32, 128, 256, 512, 1024]$  for our experiments on representation sizes. For all other experiments  $f = 256$ .

**CNN Image Decoders:** We consider a family of decoders with transposed convolutions with three TranConv-BatchNorm-ReLU layers followed by a TranConv-BatchNorm-Sigmoid layer. Filter widths are  $[128, 64, 32, 3]$  and paddings are "valid". For input sizes of  $32 \times 32$  (Cifar10, Cifar100), kernel sizes are  $[4, 4, 4, 3]$  and strides are  $[2, 2, 2, 1]$ ; for input sizes of  $64 \times 64$  (3dshapes, PCam), kernel sizes are  $[4, 4, 5, 4]$  and strides are  $[2, 2, 2, 2]$ . Weights are initialized with the standard TensorFlow initialization (kernel\_initializer="glorot\_uniform", bias\_initializer="zeros").

**CNN/ResNet Head for Rotation:** For our CNN encoder we use a fully-connected layer with 4 neurons and softmax activation as head to predict the four different rotations. We use the standard TensorFlow initialization (kernel\_initializer="glorot\_uniform", bias\_initializer="zeros"). For our ResNet decoder we initialize with (kernel\_initializer=RandomNormal(stddev=.01), bias\_initializer="zeros").

**CNN/ResNet Heads for Contrastive Learning:** For our CNN encoder we use a two layer MLP with a FC-BatchNorm-ReLU layer followed by a FC-BatchNorm-Softmax layer as projection head for contrastive learning. Number of neurons are  $[f, 128]$ . We use the standard TensorFlow initialization (kernel\_initializer="glorot\_uniform", bias\_initializer="zeros"). We vary  $f$  in  $[4, 32, 128, 256, 512, 1024]$  for our experiments on representation sizes. For all other experiments  $f = 256$ . For our ResNet head the number of neurons are  $[512, 128]$  and we initialize as in [9].

**Target Models:** For our linear target model we use a fully-connected layer with *num\_classes* neurons and a softmax activation. For our two- and three-layer nonlinear models we add layers consisting of  $[256]$  and  $[512, 256]$  hidden units with batch normalization followed by ReLU activations respectively. Weights are initialized with the standard TensorFlow initialization (kernel\_initializer="glorot\_uniform", bias\_initializer="zeros").

**Hardware:** We carry out our experiments on two servers which contain four Nvidia GeForce RTX 2080 Ti GPUs respectively.

**Mismatch Evaluation:** In Table 4 we show additional details about training, evaluation and measurements.

**Table 4** Information about measurements and training

	Measurement Epochs	Convergence Criterion	Pretext Model Training Epochs	Target Model Training Epochs	Validation
<b>Rep. Size, TMC, Augs</b>					
CAE(Cifar10)	(0,5,20,50,100,...400)	Patience:3	400	500	5-fold cross-validation
DCAE(Cifar10)	(0,5,20,50,100,...400)	Patience:3	400	500	5-fold cross-validation
CCAE(Cifar100)	(0,5,20,50,100,...400)	Patience:6	400	500	5-fold cross-validation
CCAE(PCam)	(0,10,50,100,150,200,300,...800)	Patience:10	800	500	5-fold cross-validation
RCAE(PCam)	(0,200,400,...2000)	Patience:30	2000	500	5-fold cross-validation
SCLCAE(3dshapes)	(0,10,50,100,150,200,300,...600)	Patience:15	600	100	5-fold cross-validation
<b>Target Task Type</b>					
CAE(3dshapes)	(0,10,30,50,100,...400)	Patience:3	400	100	5-fold cross-validation
DCAE(3dshapes)	(0,10,30,50,100,...400)	Patience:3	400	100	5-fold cross-validation
CCAE(3dshapes)	(0,10,30,50,100,...400)	Patience:3	400	100	5-fold cross-validation
RCAE(3dshapes)	(0,10,30,50,100,...400)	Patience:3	400	100	5-fold cross-validation
SCLCAE(3dshapes)	(0,10,30,50,100,...600)	Patience:3	600	100	5-fold cross-validation
<b>Stability</b>					
CAE100E(Cifar10)	(0,1,...,100)	Epoch 100	100	500	5-fold cross-validation
CAE(Cifar10)	(0,5,20,50,100,...400)	Epoch 400	400	500	5-fold cross-validation
CAENoCrossVal(Cifar10)	(0,5,20,50,100,...400)	Epoch 400	400	500	5 × same split
<b>ResNets</b>					
RResNet18(Cifar10)	(0,50,100,200,400,...4000)	Patience:30	4000	600	5-fold cross-validation
SCLResNet18(Cifar10)	(0,50,100,200,400,...4000)	Epoch 4000	4000	600	3-fold cross-validation
SCLResNet18(Cifar100)	(0,50,100,200,400,...4000)	Epoch 4000	4000	600	3-fold cross-validation
SCLResNet18(PCam)	(0,400,800,...5000)	Patience:60	5000	500	5-fold cross-validation
<b>ResNets Rep. Size</b>					
RResNet18(Cifar10)	(0,100,200,...,1000,1200,...3000)	Patience:30	3000	700	5-fold cross-validation

Under (Rep. Size, TMC, Augs) we refer to all models trained with different representations, target model complexities and augmentations.

## C Additional Evidence

Table 5 Detailed version of CAE (Cifar10) and DCAE (Cifar10) from Table 1

	CAE (Cifar10)			DCAE (Cifar10)		
	ACC	cSM3	MOFM	ACC	cSM3	MOFM
<b>Rep. Size</b>						
2x2x4	27.94 <sup>+1.47</sup> <sub>-0.76</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+29.77</sup> <sub>-0.00</sub>	28.06 <sup>+0.70</sup> <sub>-1.65</sub>	0.05 <sup>+0.12</sup> <sub>-0.05</sub>	<b>0.00</b> <sup>+14.27</sup> <sub>-0.00</sub>
2x2x32	36.57 <sup>+0.92</sup> <sub>-0.92</sub>	0.07 <sup>+0.09</sup> <sub>-0.07</sub>	1.99 <sup>+3.54</sup> <sub>-1.44</sub>	36.20 <sup>+0.52</sup> <sub>-0.50</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	3.20 <sup>+10.24</sup> <sub>-0.59</sub>
2x2x128	41.94 <sup>+0.45</sup> <sub>-0.47</sub>	0.20 <sup>+0.42</sup> <sub>-0.20</sub>	10.10 <sup>+6.24</sup> <sub>-3.51</sub>	41.79 <sup>+0.46</sup> <sub>-0.49</sub>	0.06 <sup>+0.11</sup> <sub>-0.06</sub>	5.51 <sup>+6.79</sup> <sub>-5.32</sub>
2x2x256	44.69 <sup>+0.69</sup> <sub>-0.41</sub>	0.75 <sup>+0.38</sup> <sub>-0.22</sub>	11.14 <sup>+5.32</sup> <sub>-2.72</sub>	45.42 <sup>+0.55</sup> <sub>-0.62</sub>	0.69 <sup>+0.43</sup> <sub>-0.69</sub>	5.17 <sup>+2.98</sup> <sub>-2.16</sub>
2x2x512	48.13 <sup>+0.58</sup> <sub>-0.88</sub>	0.43 <sup>+0.38</sup> <sub>-0.43</sub>	5.28 <sup>+1.95</sup> <sub>-2.98</sub>	49.04 <sup>+0.29</sup> <sub>-0.29</sub>	0.36 <sup>+0.75</sup> <sub>-0.36</sub>	1.25 <sup>+1.81</sup> <sub>-1.19</sub>
2x2x1024	<b>51.42</b> <sup>+0.69</sup> <sub>-0.38</sub>	<b>0.24</b> <sup>+0.32</sup> <sub>-0.24</sub>	<b>0.25</b> <sup>+1.26</sup> <sub>-0.17</sub>	<b>53.82</b> <sup>+0.57</sup> <sub>-0.82</sub>	<b>0.03</b> <sup>+0.06</sup> <sub>-0.03</sub>	<b>0.00</b> <sup>+0.04</sup> <sub>-0.00</sub>
<b>Target Model</b>						
FC	44.69 <sup>+0.69</sup> <sub>-0.41</sub>	0.75 <sup>+0.38</sup> <sub>-0.22</sub>	11.14 <sup>+5.32</sup> <sub>-2.72</sub>	45.42 <sup>+0.55</sup> <sub>-0.62</sub>	0.69 <sup>+0.43</sup> <sub>-0.69</sub>	5.17 <sup>+2.98</sup> <sub>-2.16</sub>
2FC	56.72 <sup>+0.50</sup> <sub>-0.42</sub>	<b>0.03</b> <sup>+0.10</sup> <sub>-0.03</sub>	5.68 <sup>+2.87</sup> <sub>-3.04</sub>	57.26 <sup>+0.68</sup> <sub>-0.52</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	5.14 <sup>+2.68</sup> <sub>-1.98</sub>
3FC	<b>63.05</b> <sup>+0.74</sup> <sub>-0.72</sub>	<b>0.03</b> <sup>+0.11</sup> <sub>-0.03</sub>	<b>3.94</b> <sup>+0.93</sup> <sub>-1.61</sub>	<b>63.33</b> <sup>+0.23</sup> <sub>-0.39</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>3.17</b> <sup>+1.51</sup> <sub>-0.89</sub>
<b>Augmentations</b>						
All	44.69 <sup>+0.69</sup> <sub>-0.41</sub>	<b>0.75</b> <sup>+0.38</sup> <sub>-0.13</sub>	11.14 <sup>+5.32</sup> <sub>-0.98</sub>	45.42 <sup>+0.55</sup> <sub>-0.62</sub>	0.69 <sup>+0.43</sup> <sub>-0.69</sub>	5.17 <sup>+2.98</sup> <sub>-2.16</sub>
NoJitter	46.30 <sup>+0.50</sup> <sub>-0.46</sub>	0.99 <sup>+0.59</sup> <sub>-0.47</sub>	<b>10.99</b> <sup>+1.24</sup> <sub>-3.25</sub>	47.27 <sup>+0.46</sup> <sub>-0.50</sub>	<b>0.50</b> <sup>+0.45</sup> <sub>-0.50</sub>	2.33 <sup>+1.74</sup> <sub>-1.57</sub>
NoJitterNoFlip	<b>46.48</b> <sup>+0.51</sup> <sub>-0.75</sub>	1.00 <sup>+0.48</sup> <sub>-0.60</sub>	12.51 <sup>+1.46</sup> <sub>-1.74</sub>	<b>47.38</b> <sup>+0.36</sup> <sub>-0.44</sub>	0.55 <sup>+0.51</sup> <sub>-0.55</sub>	<b>1.73</b> <sup>+4.43</sup> <sub>-1.36</sub>

Table 6 Detailed version of CCAE (Cifar100) and RCAE (PCam) from Table 1

	CCAЕ (Cifar100)			RCAE (PCam)			
	ACC	cSM3	MOFM	ACC	cSM3	MOFM	MM3
<b>Rep. Size</b>							
2x2x4	9.66 <sup>+0.30</sup> <sub>-0.45</sub>	0.28 <sup>+0.19</sup> <sub>-0.13</sub>	1.54 <sup>+2.09</sup> <sub>-0.98</sub>	67.62 <sup>+1.18</sup> <sub>-2.05</sub>	5.38 <sup>+2.15</sup> <sub>-1.66</sub>	4.15 <sup>+37.30</sup> <sub>-4.15</sub>	-22.26 <sup>+2.96</sup> <sub>-3.02</sub>
2x2x32	17.63 <sup>+0.36</sup> <sub>-0.51</sub>	0.65 <sup>+0.28</sup> <sub>-0.64</sub>	3.64 <sup>+4.14</sup> <sub>-2.04</sub>	72.83 <sup>+1.17</sup> <sub>-1.72</sub>	3.34 <sup>+0.74</sup> <sub>-0.52</sub>	7.92 <sup>+18.08</sup> <sub>-1.46</sub>	-21.09 <sup>+1.90</sup> <sub>-2.48</sub>
2x2x128	24.36 <sup>+0.56</sup> <sub>-0.47</sub>	0.51 <sup>+0.24</sup> <sub>-0.24</sub>	0.81 <sup>+1.60</sup> <sub>-0.45</sub>	78.17 <sup>+0.63</sup> <sub>-0.80</sub>	1.03 <sup>+0.88</sup> <sub>-0.96</sub>	4.04 <sup>+4.35</sup> <sub>-0.70</sub>	-23.47 <sup>+0.55</sup> <sub>-0.65</sub>
2x2x256	28.36 <sup>+0.78</sup> <sub>-0.68</sub>	0.17 <sup>+0.37</sup> <sub>-0.17</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	79.55 <sup>+1.18</sup> <sub>-1.07</sub>	0.44 <sup>+0.79</sup> <sub>-0.44</sub>	<b>0.00</b> <sup>+3.95</sup> <sub>-0.85</sub>	-27.60 <sup>+1.22</sup> <sub>-1.00</sub>
2x2x512	32.02 <sup>+0.51</sup> <sub>-0.74</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	80.80 <sup>+0.66</sup> <sub>-0.66</sub>	0.18 <sup>+0.15</sup> <sub>-0.18</sub>	<b>0.00</b> <sup>+1.13</sup> <sub>-0.00</sub>	-28.03 <sup>+1.28</sup> <sub>-1.35</sub>
2x2x1024	<b>34.89</b> <sup>+0.41</sup> <sub>-0.83</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>82.20</b> <sup>+0.76</sup> <sub>-0.68</sub>	<b>0.09</b> <sup>+0.12</sup> <sub>-0.09</sub>	<b>0.00</b> <sup>+0.34</sup> <sub>-0.00</sub>	-26.56 <sup>+1.14</sup> <sub>-1.18</sub>
<b>Target Model</b>							
FC	32.02 <sup>+0.51</sup> <sub>-0.74</sub>	<b>0.00</b> <sup>+0.02</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	78.17 <sup>+0.63</sup> <sub>-0.80</sub>	1.03 <sup>+0.88</sup> <sub>-0.35</sub>	4.04 <sup>+4.35</sup> <sub>-0.70</sub>	-23.47 <sup>+0.55</sup> <sub>-0.65</sub>
2FC	36.01 <sup>+0.31</sup> <sub>-0.27</sub>	0.08 <sup>+0.32</sup> <sub>-0.08</sub>	<b>0.00</b> <sup>+0.39</sup> <sub>-0.00</sub>	82.99 <sup>+0.65</sup> <sub>-0.62</sub>	<b>0.31</b> <sup>+0.31</sup> <sub>-0.31</sub>	0.76 <sup>+1.76</sup> <sub>-0.04</sub>	-28.56 <sup>+0.81</sup> <sub>-0.72</sub>
3FC	<b>38.40</b> <sup>+0.34</sup> <sub>-0.50</sub>	0.12 <sup>+0.28</sup> <sub>-0.12</sub>	0.02 <sup>+0.44</sup> <sub>-0.09</sub>	<b>84.18</b> <sup>+0.62</sup> <sub>-0.45</sub>	0.37 <sup>+0.21</sup> <sub>-0.37</sub>	<b>0.61</b> <sup>+3.68</sup> <sub>-0.51</sub>	-29.61 <sup>+0.71</sup> <sub>-0.61</sub>
<b>Augmentations</b>							
All	28.36 <sup>+0.78</sup> <sub>-0.68</sub>	<b>0.17</b> <sup>+0.37</sup> <sub>-0.17</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	78.17 <sup>+0.63</sup> <sub>-0.80</sub>	1.03 <sup>+0.88</sup> <sub>-0.96</sub>	4.04 <sup>+4.35</sup> <sub>-0.70</sub>	-23.47 <sup>+0.55</sup> <sub>-0.65</sub>
NoJitter	-	-	-	<b>80.88</b> <sup>+0.88</sup> <sub>-0.94</sub>	<b>0.58</b> <sup>+0.34</sup> <sub>-0.58</sub>	<b>0.10</b> <sup>+5.52</sup> <sub>-0.04</sub>	-28.60 <sup>+2.23</sup> <sub>-1.18</sub>
NoJitterNoFlip	-	-	-	80.55 <sup>+0.73</sup> <sub>-0.55</sub>	1.51 <sup>+0.57</sup> <sub>-1.09</sub>	7.33 <sup>+9.76</sup> <sub>-1.78</sub>	-10.60 <sup>+1.21</sup> <sub>-1.88</sub>
NoFlip	<b>29.33</b> <sup>+0.65</sup> <sub>-0.61</sub>	0.20 <sup>+0.43</sup> <sub>-0.20</sub>	<b>0.00</b> <sup>+0.22</sup> <sub>-0.00</sub>	-	-	-	-

Table 7 Detailed version of CCAE (PCam) and SCLCAE (3dshapes for object hue classification) from Table 1

	CCAЕ (PCam)			SCLCAE (3dshapes)			
	ACC	cSM3	MOFM	ACC	cSM3	MOFM	MM3
<b>Rep. Size</b>							
2x2x4	63.39 <sup>+3.80</sup> <sub>-1.58</sub>	4.98 <sup>+5.53</sup> <sub>-3.53</sub>	9.28 <sup>+∞</sup> <sub>-2.09</sub>	38.07 <sup>+11.18</sup> <sub>-10.06</sub>	26.34 <sup>+11.38</sup> <sub>-9.69</sub>	∞	-7.99 <sup>+1.83</sup> <sub>-1.28</sub>
2x2x32	72.73 <sup>+2.98</sup> <sub>-2.37</sub>	5.17 <sup>+2.95</sup> <sub>-2.80</sub>	34.30 <sup>+16.31</sup> <sub>-10.21</sub>	85.25 <sup>+2.77</sup> <sub>-2.36</sub>	12.98 <sup>+10.56</sup> <sub>-12.24</sub>	36.39 <sup>+26.87</sup> <sub>-29.73</sub>	-57.67 <sup>+6.98</sup> <sub>-6.46</sub>
2x2x128	78.66 <sup>+0.46</sup> <sub>-0.22</sub>	0.32 <sup>+0.47</sup> <sub>-0.32</sub>	0.10 <sup>+3.51</sup> <sub>-0.48</sub>	96.54 <sup>+0.72</sup> <sub>-0.83</sub>	8.14 <sup>+1.29</sup> <sub>-2.53</sub>	<b>22.65</b> <sup>+13.03</sup> <sub>-3.92</sub>	-66.19 <sup>+1.24</sup> <sub>-1.71</sub>
2x2x256	79.97 <sup>+0.68</sup> <sub>-0.51</sub>	0.43 <sup>+0.43</sup> <sub>-0.43</sub>	0.87 <sup>+4.56</sup> <sub>-0.41</sub>	98.65 <sup>+0.83</sup> <sub>-0.80</sub>	6.52 <sup>+0.62</sup> <sub>-0.84</sub>	27.65 <sup>+5.66</sup> <sub>-4.87</sub>	-65.77 <sup>+2.19</sup> <sub>-1.12</sub>
2x2x512	82.34 <sup>+0.66</sup> <sub>-0.87</sub>	0.20 <sup>+0.59</sup> <sub>-0.20</sub>	0.07 <sup>+0.24</sup> <sub>-0.07</sub>	99.41 <sup>+0.20</sup> <sub>-0.25</sub>	5.96 <sup>+1.13</sup> <sub>-0.40</sub>	27.78 <sup>+8.80</sup> <sub>-7.10</sub>	-63.61 <sup>+1.75</sup> <sub>-1.50</sub>
2x2x1024	<b>83.67</b> <sup>+0.45</sup> <sub>-0.61</sub>	<b>0.00</b> <sup>+0.01</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.20</sup> <sub>-0.00</sub>	<b>99.75</b> <sup>+0.04</sup> <sub>-0.04</sub>	<b>4.76</b> <sup>+1.84</sup> <sub>-1.36</sub>	32.70 <sup>+14.59</sup> <sub>-7.15</sub>	-57.92 <sup>+2.29</sup> <sub>-2.25</sub>
<b>Target Model</b>							
FC	83.67 <sup>+0.45</sup> <sub>-0.61</sub>	<b>0.00</b> <sup>+0.01</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.20</sup> <sub>-0.00</sub>	38.07 <sup>+11.18</sup> <sub>-10.06</sub>	6.52 <sup>+0.62</sup> <sub>-0.84</sub>	<b>27.65</b> <sup>+5.66</sup> <sub>-4.87</sub>	-65.77 <sup>+2.19</sup> <sub>-1.12</sub>
2FC	89.17 <sup>+0.50</sup> <sub>-0.38</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	85.25 <sup>+2.77</sup> <sub>-2.36</sub>	1.84 <sup>+0.33</sup> <sub>-0.66</sub>	103.30 <sup>+36.89</sup> <sub>-21.22</sub>	-70.49 <sup>+1.41</sup> <sub>-0.75</sub>
3FC	<b>90.59</b> <sup>+0.48</sup> <sub>-0.60</sub>	0.08 <sup>+0.19</sup> <sub>-0.08</sub>	<b>0.00</b> <sup>+0.44</sup> <sub>-0.00</sub>	<b>96.54</b> <sup>+0.72</sup> <sub>-0.83</sub>	<b>0.91</b> <sup>+0.18</sup> <sub>-0.30</sub>	258.18 <sup>+308.48</sup> <sub>-88.99</sub>	-71.26 <sup>+1.20</sup> <sub>-0.77</sub>
<b>Augmentations</b>							
All	79.97 <sup>+0.68</sup> <sub>-0.51</sub>	0.43 <sup>+0.51</sup> <sub>-0.43</sub>	0.87 <sup>+4.56</sup> <sub>-0.41</sub>	38.07 <sup>+11.18</sup> <sub>-10.06</sub>	6.52 <sup>+0.62</sup> <sub>-0.84</sub>	27.65 <sup>+5.66</sup> <sub>-4.87</sub>	-65.77 <sup>+2.19</sup> <sub>-1.12</sub>
NoJitter	-	-	-	85.25 <sup>+2.77</sup> <sub>-2.36</sub>	<b>0.00</b> <sup>+0.02</sup> <sub>-0.00</sub>	0.01 <sup>+0.07</sup> <sub>-0.01</sub>	-37.92 <sup>+0.72</sup> <sub>-0.89</sub>
NoJitterNoFlip	-	-	-	<b>96.54</b> <sup>+0.72</sup> <sub>-0.83</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.04</sup> <sub>-0.01</sub>	-35.95 <sup>+0.63</sup> <sub>-0.82</sub>
NoFlip	<b>80.70</b> <sup>+0.36</sup> <sub>-0.57</sub>	<b>0.41</b> <sup>+0.21</sup> <sub>-0.30</sub>	<b>0.04</b> <sup>+1.03</sup> <sub>-0.20</sub>	-	-	-	-

**Table 8** Detailed version of CAE and DCAE from Table 2

	ACC	CAE cSM3	MOFM	ACC	DCAE cSM3	MOFM
floor_hue	99.96 <sup>+0.02</sup> <sub>-0.02</sub>	0.01 <sup>+0.02</sup> <sub>-0.01</sub>	0.95 <sup>+0.71</sup> <sub>-0.51</sub>	99.97 <sup>+0.02</sup> <sub>-0.06</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	1.28 <sup>+1.11</sup> <sub>-0.81</sub>
wall_hue	<b>100.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	0.02 <sup>+0.04</sup> <sub>-0.02</sub>	32.03 <sup>+6.38</sup> <sub>-8.59</sub>	<b>99.99</b> <sup>+0.00</sup> <sub>-0.01</sub>	<b>0.00</b> <sup>+0.01</sup> <sub>-0.00</sub>	24.43 <sup>+7.56</sup> <sub>-6.97</sub>
object_hue	99.23 <sup>+0.18</sup> <sub>-0.35</sub>	0.38 <sup>+0.82</sup> <sub>-0.37</sub>	22.71 <sup>+14.77</sup> <sub>-6.45</sub>	99.22 <sup>+0.26</sup> <sub>-0.31</sub>	0.43 <sup>+0.61</sup> <sub>-0.43</sub>	24.55 <sup>+7.04</sup> <sub>-11.17</sub>
scale	74.17 <sup>+6.07</sup> <sub>-6.42</sub>	0.41 <sup>+1.64</sup> <sub>-0.41</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	68.27 <sup>+2.36</sup> <sub>-2.04</sub>	0.27 <sup>+0.59</sup> <sub>-0.27</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>
shape	98.38 <sup>+0.82</sup> <sub>-1.19</sub>	0.07 <sup>+0.14</sup> <sub>-0.07</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	97.45 <sup>+0.46</sup> <sub>-0.57</sub>	0.08 <sup>+0.20</sup> <sub>-0.08</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>
orientation	81.63 <sup>+3.23</sup> <sub>-2.89</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	74.27 <sup>+1.07</sup> <sub>-2.34</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>
average	<b>92.22</b>	0.15	9.28	89.86	<b>0.13</b>	8.21

**Table 9** Detailed version of CCAE and RCAE from Table 2

	ACC	CCA cSM3	MOFM	ACC	RCA cSM3	MOFM	MM3
floor_hue	99.94 <sup>+0.03</sup> <sub>-0.05</sub>	<b>0.02</b> <sup>+0.03</sup> <sub>-0.02</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	92.08 <sup>+1.79</sup> <sub>-2.89</sub>	56.68 <sup>+0.94</sup> <sub>-1.55</sub>	∞	44.67 <sup>+4.10</sup> <sub>-3.78</sub>
wall_hue	<b>99.98</b> <sup>+0.01</sup> <sub>-0.02</sub>	0.10 <sup>+0.11</sup> <sub>-0.07</sub>	<b>0.00</b> <sup>+0.78</sup> <sub>-0.00</sub>	<b>99.96</b> <sup>+0.04</sup> <sub>-0.04</sub>	25.17 <sup>+9.67</sup> <sub>-8.15</sub>	∞	7.80 <sup>+0.59</sup> <sub>-0.85</sub>
object_hue	98.96 <sup>+0.38</sup> <sub>-0.30</sub>	1.55 <sup>+0.73</sup> <sub>-0.63</sub>	0.63 <sup>+5.38</sup> <sub>-0.59</sub>	98.79 <sup>+0.35</sup> <sub>-0.41</sub>	59.65 <sup>+8.12</sup> <sub>-10.39</sub>	∞	40.11 <sup>+0.99</sup> <sub>-1.88</sub>
scale	66.72 <sup>+3.20</sup> <sub>-2.20</sub>	0.10 <sup>+0.41</sup> <sub>-0.10</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	67.10 <sup>+6.77</sup> <sub>-3.96</sub>	2.60 <sup>+4.10</sup> <sub>-2.60</sub>	0.13 <sup>+3.54</sup> <sub>-0.13</sub>	31.78 <sup>+1.91</sup> <sub>-2.44</sub>
shape	98.75 <sup>+0.39</sup> <sub>-0.50</sub>	0.03 <sup>+0.10</sup> <sub>-0.03</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	98.17 <sup>+0.43</sup> <sub>-0.39</sub>	<b>0.20</b> <sup>+0.45</sup> <sub>-0.20</sub>	0.06 <sup>+0.82</sup> <sub>-0.06</sub>	-2.48 <sup>+0.65</sup> <sub>-1.28</sub>
orientation	72.55 <sup>+3.17</sup> <sub>-2.81</sub>	0.23 <sup>+0.91</sup> <sub>-0.23</sub>	<b>0.00</b> <sup>+0.00</sup> <sub>-0.00</sub>	80.69 <sup>+2.35</sup> <sub>-2.68</sub>	0.48 <sup>+1.91</sup> <sub>-0.48</sub>	<b>0.00</b> <sup>+1.23</sup> <sub>-0.00</sub>	22.26 <sup>+0.74</sup> <sub>-1.19</sub>
average	89.48	0.34	<b>0.11</b>	89.47	24.13	∞	24.02

**Table 10** Detailed version of SCLCAE from Table 2

	ACC	cSM3	MOFM	MM3
floor_hue	93.53 <sup>+2.62</sup> <sub>-1.90</sub>	28.18 <sup>+11.12</sup> <sub>-6.30</sub>	268.27 <sup>+161.92</sup> <sub>-44.43</sub>	-48.38 <sup>+2.75</sup> <sub>-2.15</sub>
wall_hue	<b>99.96</b> <sup>+0.03</sup> <sub>-0.03</sub>	<b>0.29</b> <sup>+0.07</sup> <sub>-0.11</sub>	0.46 <sup>+5.04</sup> <sub>-0.44</sub>	<b>-76.40</b> <sup>+2.97</sup> <sub>-3.06</sub>
object_hue	98.67 <sup>+0.82</sup> <sub>-1.06</sub>	2.87 <sup>+0.69</sup> <sub>-1.36</sub>	8.69 <sup>+5.20</sup> <sub>-0.64</sub>	-73.17 <sup>+2.58</sup> <sub>-2.77</sub>
scale	83.94 <sup>+1.15</sup> <sub>-0.57</sub>	2.43 <sup>+4.02</sup> <sub>-2.43</sub>	<b>0.00</b> <sup>+1.78</sup> <sub>-0.00</sub>	-44.80 <sup>+2.06</sup> <sub>-1.55</sub>
shape	95.06 <sup>+0.92</sup> <sub>-0.86</sub>	1.67 <sup>+1.54</sup> <sub>-1.67</sub>	2.16 <sup>+1.48</sup> <sub>-0.13</sub>	-67.54 <sup>+1.38</sup> <sub>-1.63</sub>
orientation	45.32 <sup>+4.14</sup> <sub>-3.57</sub>	2.50 <sup>+2.53</sup> <sub>-1.97</sub>	6.68 <sup>+3.21</sup> <sub>-3.02</sub>	-9.11 <sup>+2.46</sup> <sub>-1.71</sub>
average	86.08	6.32	47.71	<b>-53.23</b>

**Table 11** Mismatches of other models we have tested

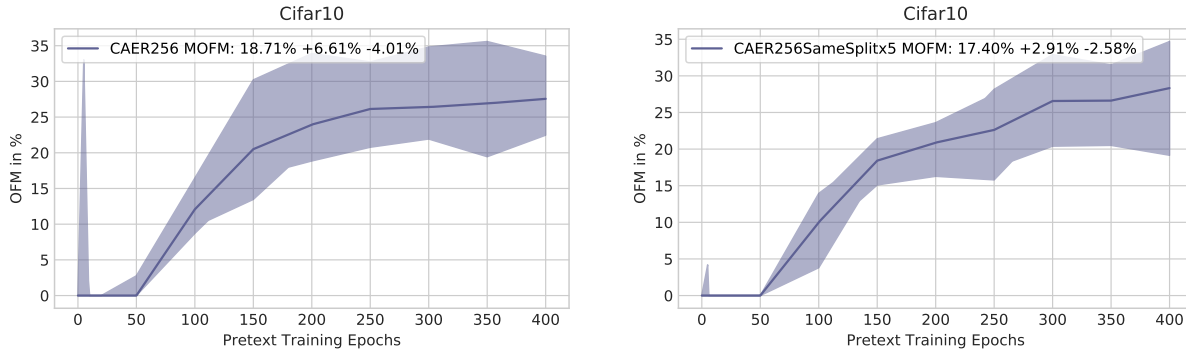
	Cifar10			
	ACC	cSM3	MOFM	MM3
CAE	44.69 <sup>+0.69</sup> <sub>-0.41</sub>	1.42 <sup>+0.68</sup> <sub>-0.70</sub>	18.71 <sup>+6.61</sup> <sub>-4.01</sub>	-
CAENoCrossVal	44.24 <sup>+0.24</sup> <sub>-0.41</sub>	1.26 <sup>+0.15</sup> <sub>-0.47</sub>	17.40 <sup>+2.91</sup> <sub>-2.58</sub>	-
CAE100E	45.37 <sup>+1.09</sup> <sub>-1.06</sub>	1.22 <sup>+0.56</sup> <sub>-0.32</sub>	5.84 <sup>+7.59</sup> <sub>-1.12</sub>	-
RResNet18	54.64 <sup>+1.80</sup> <sub>-2.01</sub>	3.98 <sup>+1.74</sup> <sub>-3.60</sub>	4.87 <sup>+4.42</sup> <sub>-3.11</sub>	31.82 <sup>+0.75</sup> <sub>-0.68</sub>
SCLResNet18	87.14 <sup>+0.37</sup> <sub>-0.40</sub>	0.29 <sup>+0.23</sup> <sub>-0.20</sub>	0.12 <sup>+0.07</sup> <sub>-0.02</sub>	-31.83 <sup>+0.12</sup> <sub>-0.23</sub>
RResNet18R32	39.14 <sup>+1.53</sup> <sub>-1.69</sub>	1.92 <sup>+1.75</sup> <sub>-1.29</sub>	4.39 <sup>+9.53</sup> <sub>-1.58</sub>	48.93 <sup>+1.17</sup> <sub>-2.19</sub>
RResNet18R256	47.46 <sup>+1.44</sup> <sub>-1.76</sub>	3.91 <sup>+2.94</sup> <sub>-1.71</sub>	6.86 <sup>+8.51</sup> <sub>-2.09</sub>	39.08 <sup>+1.78</sup> <sub>-2.07</sub>
RResNet18R512	50.55 <sup>+2.91</sup> <sub>-2.84</sub>	7.01 <sup>+1.83</sup> <sub>-3.73</sub>	9.83 <sup>+10.55</sup> <sub>-3.90</sub>	36.93 <sup>+3.15</sup> <sub>-1.96</sub>
RResNet18R756	51.80 <sup>+2.04</sup> <sub>-2.38</sub>	7.56 <sup>+5.01</sup> <sub>-3.17</sub>	8.50 <sup>+11.90</sup> <sub>-3.00</sub>	35.61 <sup>+1.45</sup> <sub>-1.74</sub>
RResNet18R1024	52.99 <sup>+1.54</sup> <sub>-2.58</sub>	9.89 <sup>+4.51</sup> <sub>-2.41</sub>	15.12 <sup>+11.12</sup> <sub>-4.18</sub>	36.86 <sup>+2.76</sup> <sub>-2.20</sub>

Values are obtained by cross-validation, please refer to Table 4 for more details.

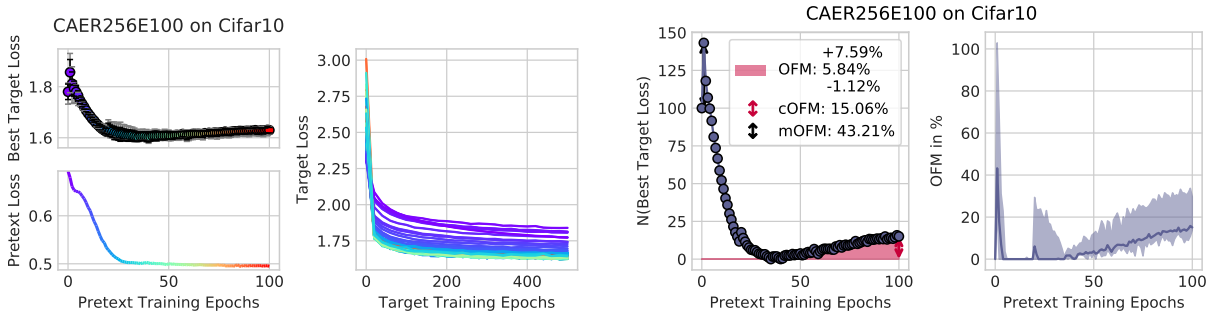
**Table 12** Mismatches of other models we have tested

	PCam				Cifar100			
	ACC	cSM3	MOFM	MM3	ACC	cSM3	MOFM	MM3
SCLResNet18	96.25 <sup>+0.44</sup> <sub>-0.23</sub>	0.37 <sup>+0.44</sup> <sub>-0.37</sub>	0.86 <sup>+1.00</sup> <sub>-0.60</sub>	-53.26 <sup>+0.52</sup> <sub>-0.38</sub>	59.20 <sup>+0.19</sup> <sub>-0.25</sub>	0.41 <sup>+0.10</sup> <sub>-0.09</sub>	0.06 <sup>+0.05</sup> <sub>-0.04</sub>	0.84 <sup>+0.06</sup> <sub>-0.03</sub>

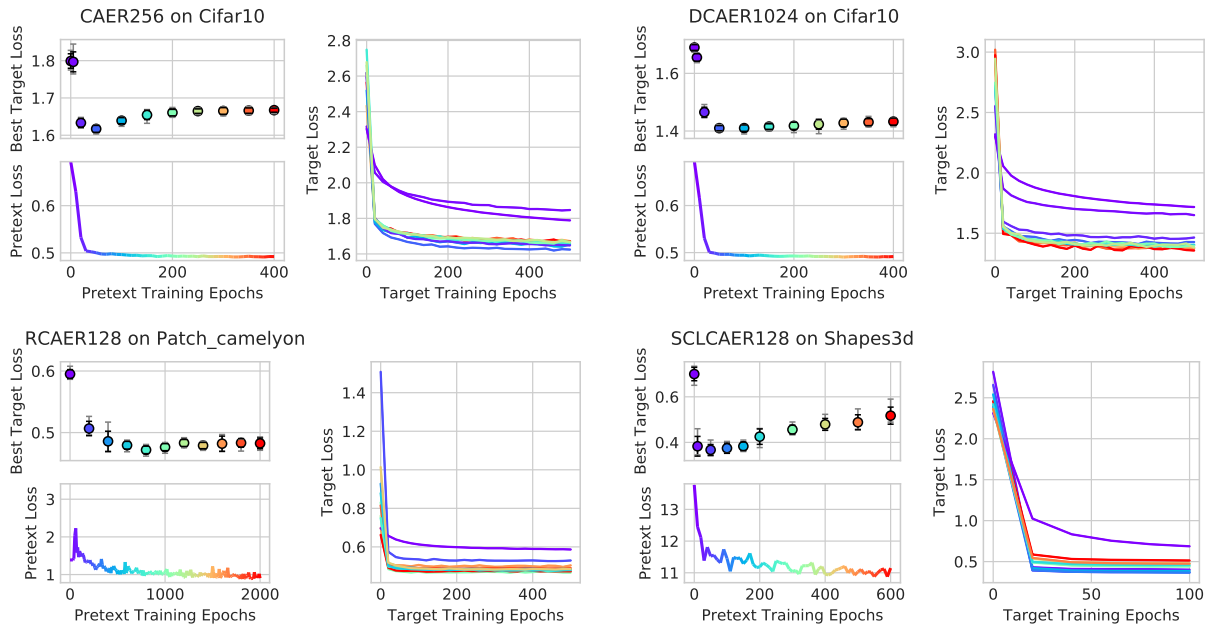
Values are obtained by cross-validation, please refer to Table 4 for more details.



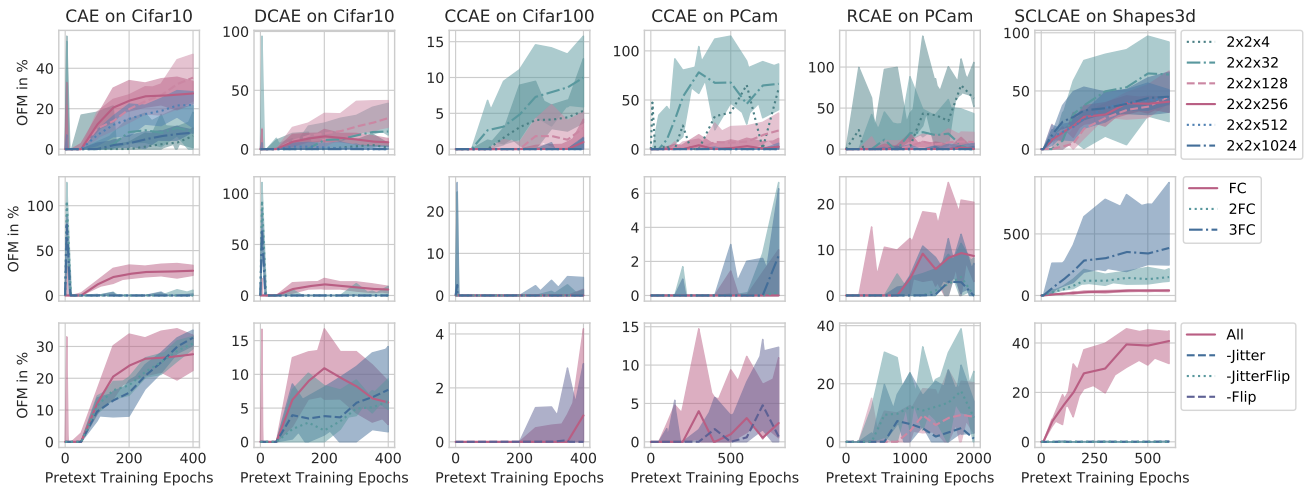
**Fig. 7** (Left) OFM of a CAE trained for 400 epochs. Stability is measured by 5-fold cross-validation. (Right) OFM of a CAE trained for 400 epochs. Stability is measured by training the CAE five times on the same dataset split. Unsurprisingly, the stability of the OFM is higher when the CAE is trained on the same split instead of the different splits from the 5-fold cross-validation.



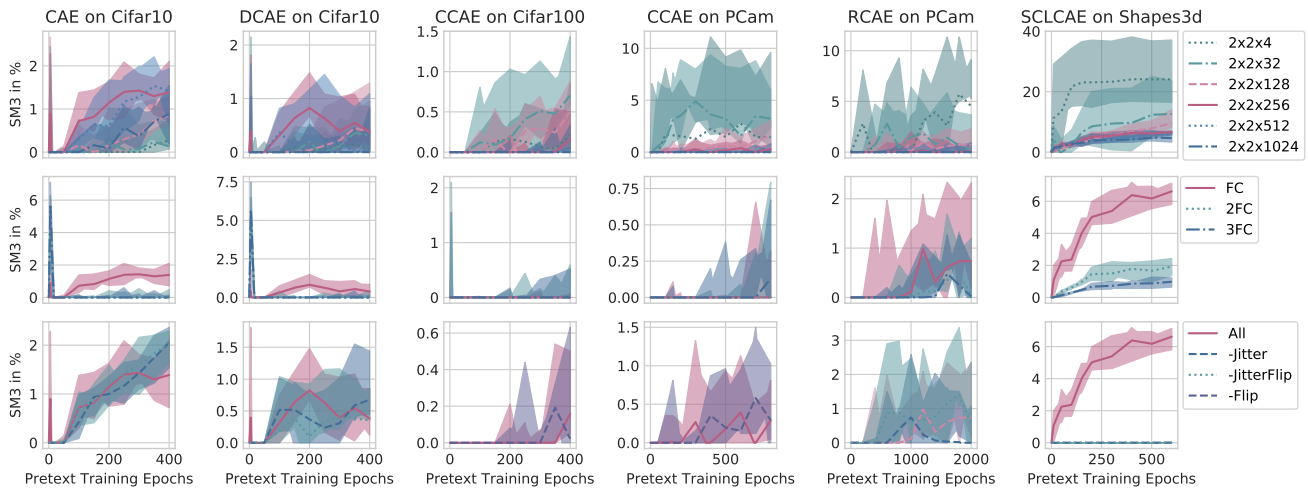
**Fig. 8** (Left) Losses of a simple CAE measured for every pretext training epoch. The curve formed by the target models represents a typical target training curve in our setup. (Right) The OFM and its stability measured for every pretext training epoch of the CAE. When we compare the stability to the partially measured CAE in Figure 7, we observe a similar instability.



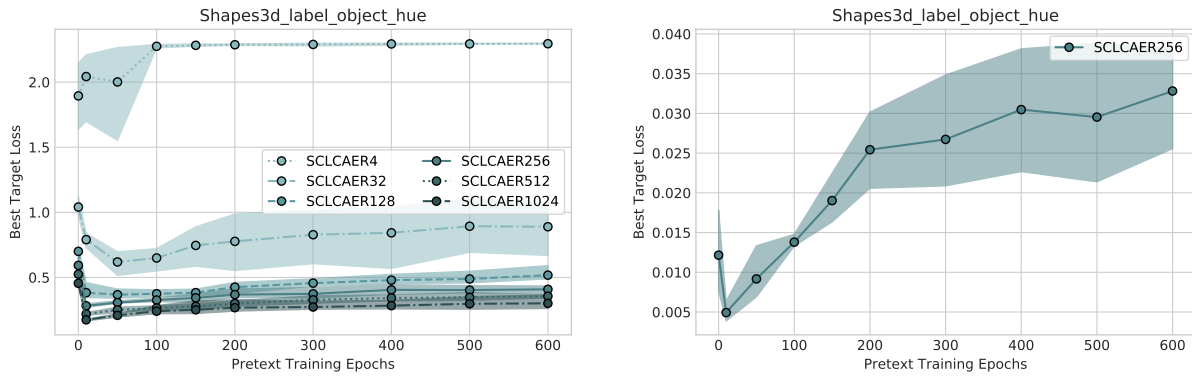
**Fig. 9** Additional evidence for the Mismatch and Convergence Section 6.1: *Longer training of the pretext task tends to create easier separable representations which may mismatch with the class label.* We observe that the target loss curves converge faster for target models trained on the pretext model’s representation from later pretext training epochs. Especially the purple curves show this behavior clearly.



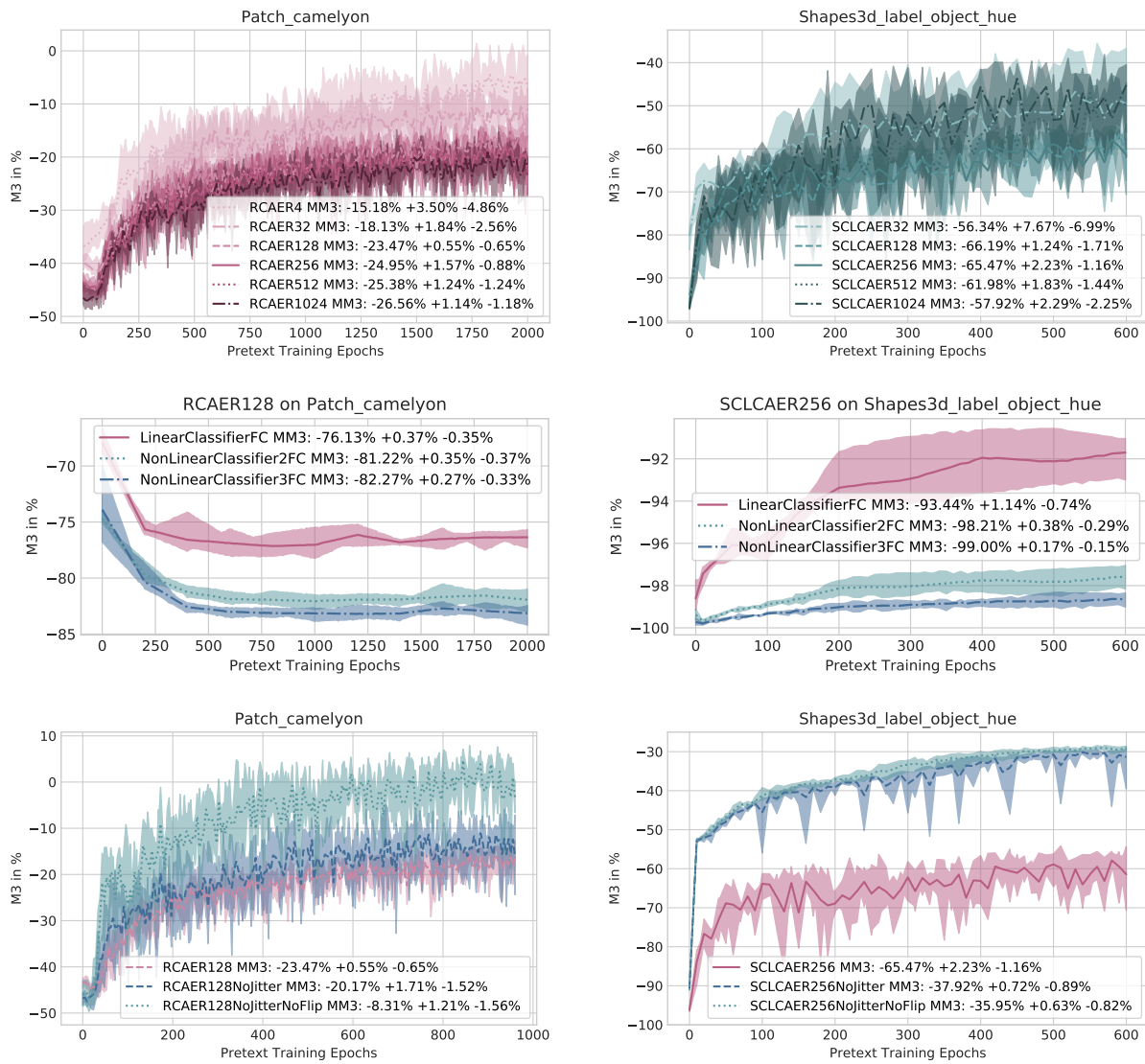
**Fig. 10** Version of Figure 4 without convergence criterium. We observe similar behaviors of the mismatches as in Figure 4 in most cases. One exception is the color jitter, where the OFM starts to converge or decrease late in training. (Top) Impact of different pretext model representation sizes on the OFM for our model. (Middle) The OFM for the linear and nonlinear target models trained on our pretext model. (Bottom) The OFM for the linear target model and for the pretext models trained on fewer augmentations. First we removed the color jitter and then the vertical flip from the augmentations. For the CCAE we only removed the vertical flip. The target models of SCLCAE were trained on 3dshapes, to predict the object hue.



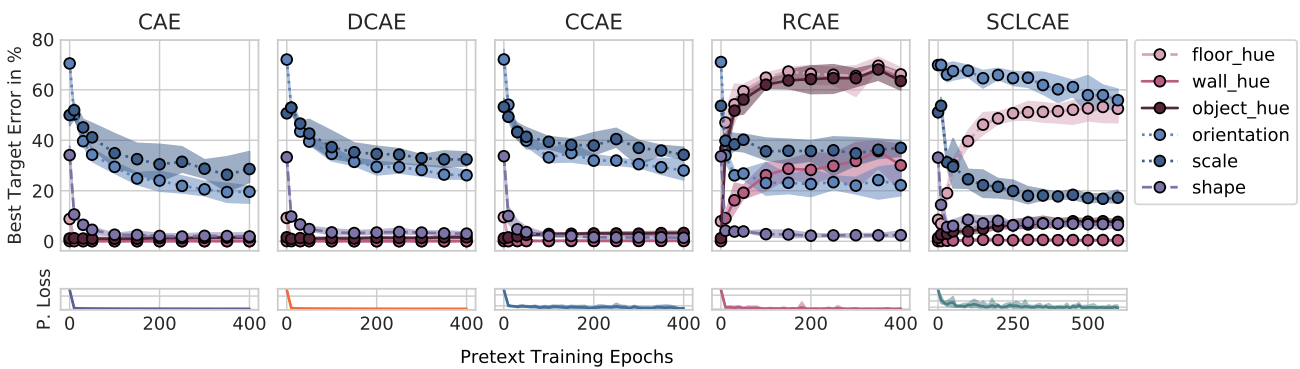
**Fig. 11** Version of Figure 4 for SM3 on accuracies, without convergence criterium. Again, we observe the similar behaviors of the mismatches as in Figure 4 in most cases.



**Fig. 12** (Left) We describe why the *OFM* for higher representation sizes does not decrease in the setup where we train the contrastive pretext model (SCLCAE) on 3dshaped and the target model to predict the object hue. Here, the target models trained on the untrained pretext models with larger representation sizes already achieve high performance due to a higher amount of color-selective, random features. Additionally, learning of the pretext model does not lead to a high performance gain, which leads again to a small interval for normalization. Therefore, forgetting useful features for the target task later in training leads to a high mismatch. (Right) We describe why the *OFM* does not decrease for more complex target models in the same setup. The nonlinear target model can make better sense of specific random pretext features for classification, which leads to a very low target loss at pretext model initialization. Since the pretext model does not learn many useful features for the target task later, this leads again to a small interval for normalization. Therefore, the *OFM* gets very large later in training, when the pretext model starts to forget useful features for the target task.

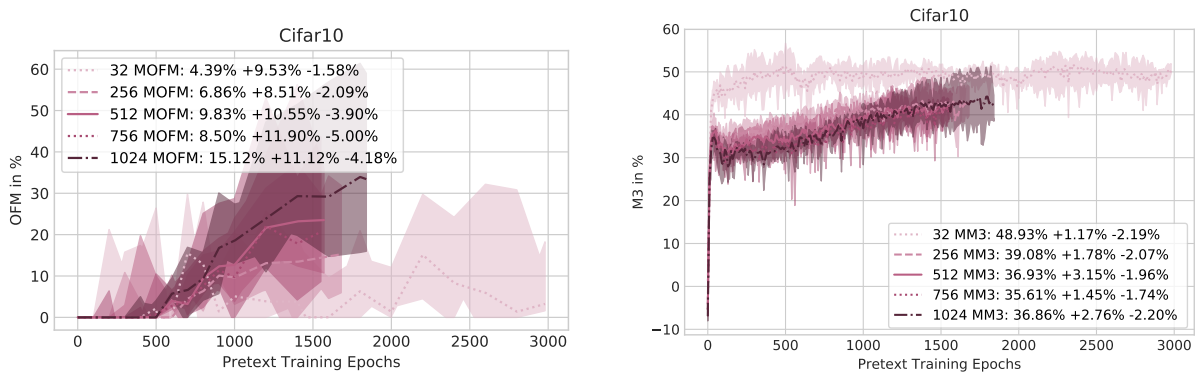


**Fig. 13** Version of Figure 4 for M3 on accuracies. We observe that besides early spikes, M3 decreases when we add complexity to the target model. Additionally, we observe that M3 decreases when we add augmentations in this case. We measure the M3 for RCAE between the classification error of the target task and the classification error of predicting the rotations of rotated images from PCam (pretext task). For SCLCAE the pretext task metric measures the ability of the model to correctly detect the representation of each given image in a batch of representations of transformed images. Here we show M3 without a convergence criterium.

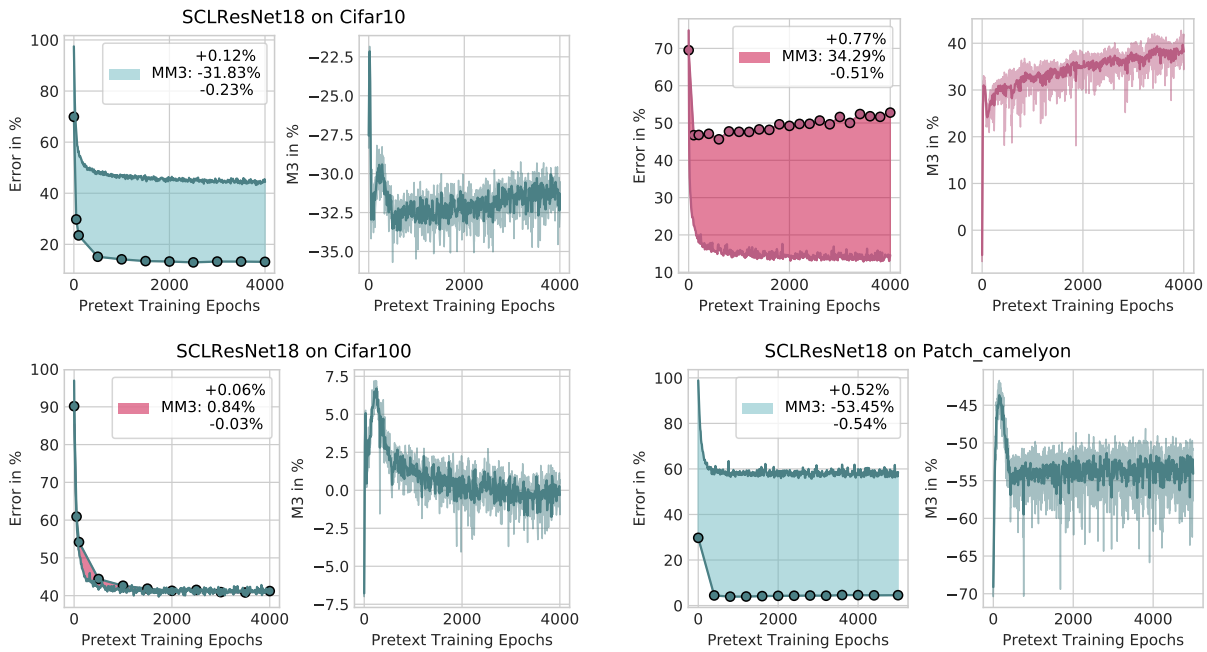


**Fig. 14** Version of Figure 5 for SM3 on accuracies.





**Fig. 15** (Left) The OFM of ResNets with different representation sizes trained on the pretext task of predicting rotations on Cifar10. Target models are trained for Cifar10 classification. (Right) MM3 of those ResNets. We vary representation sizes in [32, 256, 512, 756, 1024] by adding a  $1 \times 1$  convolution layer on top of the ResNet18. Thereby the number of filters corresponds to the representation size. In contrast to our observations on our small model task the largest representation we tested leads to a high OFM. A reason for that could be that the larger representation size helps the model to solve the pretext task and since there is a mismatch with the target task, a better understanding of this task leads to a higher mismatch. We note that a representation size of 1024 is still very small for unsupervised learning. An even larger representation size could therefore still lead to a lower mismatch.



**Fig. 16** MM3 for different pretext tasks trained with a ResNet18 model as backbone. The mismatches are shown for the entire training process. In contrast to the prediction of rotations (top right), SCLResNet18 has a high negative MM3 for Cifar10. This indicates that learning the contrastive pretext task is better suited for distinguishing Cifar10 classes than the prediction of rotations. Furthermore, the error of the contrastive pretext task is significant, which indicates that the model still underfits the pretext task with this setup and there is more room for improvement. For the 100 classes of Cifar100, MM3 becomes slightly positive in the contrastive learning setup. For contrastive learning on the PCam dataset and rotation prediction on Cifar10, we observe an increasing mismatch during training.