

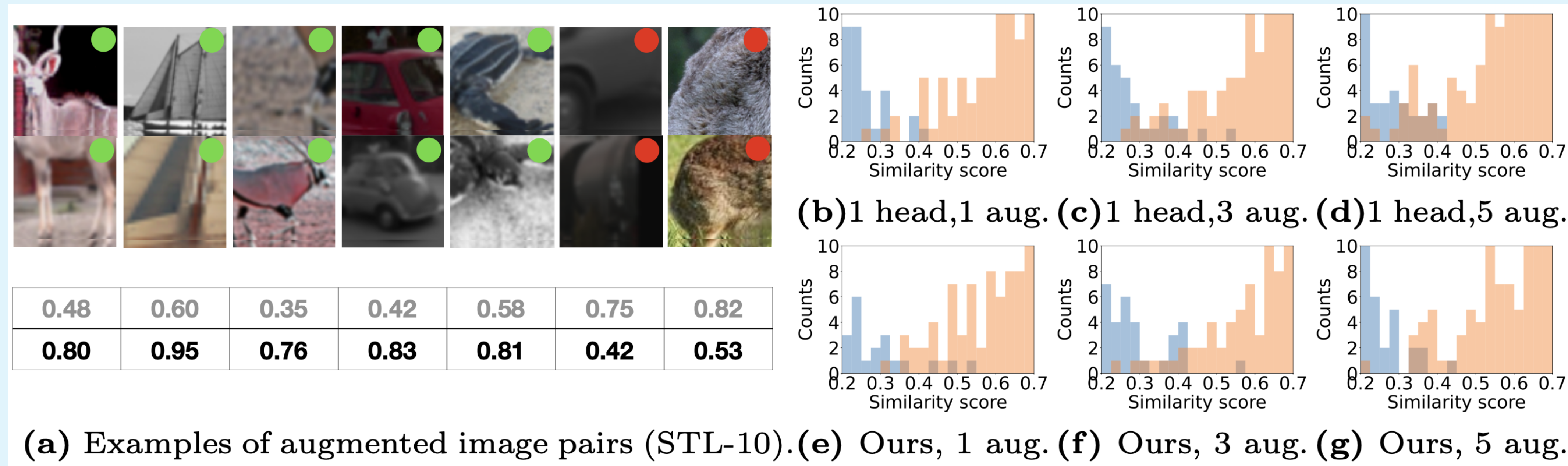
Adaptive Multi-head Contrastive Learning

Lei Wang^{1, 2} Piotr Koniusz^{2, 1} Tom Gedeon³ Liang Zheng¹

¹Australian National University ²Data61/CSIRO ³Curtin University



Motivation and key ideas



- Diverse augmentation strategies in contrastive learning and varying intra-sample similarity cause views from the same image may not always be similar.
- Owing to inter-sample similarity, views from different images may be more akin than those from the same image.
- The table in (a) shows the original (gray) and our method's similarity scores (black).
- (b)-(d): for traditional contrastive learning methods, when increasing the number of augmentations from 1 to 5, similarities of more positive pairs drop below 0.5, causing more significant overlapping regions between histograms of positive (orange) and negative (blue) pairs.
- In comparison, our multi-head approach (e)-(g) yields better separation of positive and negative sample pairs as more augmentation types are used, e.g., (g) vs (d).

Standard contrastive learning methods and their loss functions:

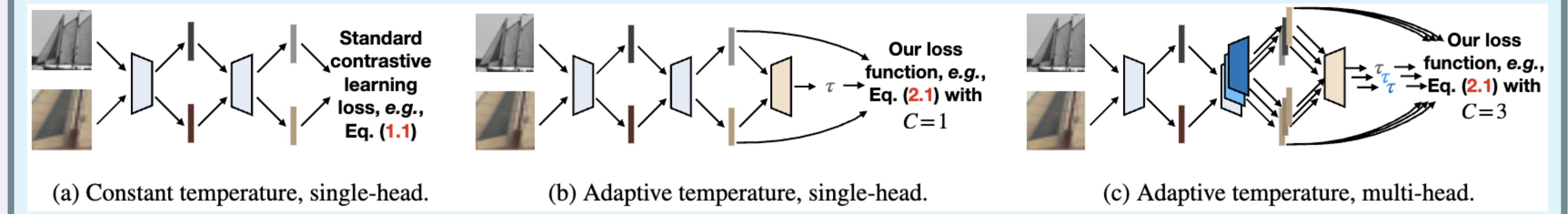
Method	Loss name	Loss function
SimCLR, MoCo	NT-Xent	$\ell_{\text{NT-Xent}} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_i^+)/\tau)}{\sum_{n=1}^N \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_{in}^-)/\tau)}$
SimSiam	Negative cos.	$\ell_{\text{SymNegCos}} = -\frac{1}{2} \text{sim}(\mathbf{z}_i, [\mathbf{h}_i^+]_{\text{sg}}) - \frac{1}{2} \text{sim}(\mathbf{z}_i^+, [\mathbf{h}_i]_{\text{sg}})$
Barlow Twins	Cross-corr.	$\ell_{\text{Cross-Corr}} = \sum_{l=1}^{d'} (1 - \mathcal{C}_{ll})^2 + \lambda \sum_{l=1}^{d'} \sum_{m \neq l}^{d'} \mathcal{C}_{lm}^2$
LGP, CAN	InfoNCE	$\ell_{\text{InfoNCE}} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_i^+)/\tau)}{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_i^+)/\tau) + \sum_{n=1}^N \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_{in}^-)/\tau)}$

- We propose adaptive multi-head contrastive learning (AMCL): it better captures the diverse image content and gives similarity scores that better separate positive and negative pairs.
- Within AMCL, we design an adaptive temperature which depends on both the projection head and the similarity of the current pair.

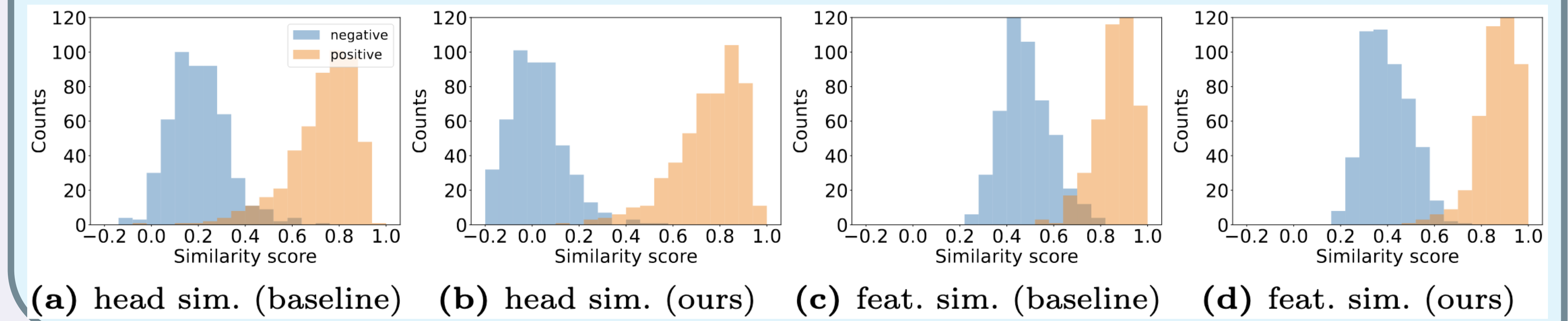
Loss functions for applying AMCL to widely used contrastive learning frameworks involve introducing C heads and a regularization term (highlighted in green):

Method	Loss function	Regularization
SimCLR, MoCo	$\ell_{\text{NT-Xent}}^\dagger = \sum_{c=1}^C \left(-\frac{1}{\tau_i^{c+}} \text{sim}(\mathbf{z}_i^c, \mathbf{z}_i^{c+}) + \frac{1}{\tau_{in*}^{c-}} \max_{n=1, \dots, N} \text{sim}(\mathbf{z}_i^c, \mathbf{z}_{in}^{c-}) \right) + \beta \Omega(\tau_i^{c+}) - \beta \Omega(\tau_{in*}^{c-})$	
SimSiam	$\ell_{\text{SymNegCos}}^\dagger = \sum_{c=1}^C \left(-\frac{1}{2\tau_i^{c+}} \text{sim}(\mathbf{z}_i^c, [\mathbf{h}_i^+]_{\text{sg}}) - \frac{1}{2\tau_{in*}^{c-}} \text{sim}(\mathbf{z}_i^{c+}, [\mathbf{h}_i]_{\text{sg}}) \right) + \beta \Omega(\tau_i^{c+}) + \beta \Omega(\tau_{in*}^{c-})$	
Barlow Twins	$\ell_{\text{Cross-Corr}}^\dagger = \sum_{c=1}^C \left(\sum_{l=1}^{d'} (1 - \frac{1}{\tau_l^{c+}} \mathcal{C}_{ll})^2 + \lambda \sum_{l=1}^{d'} \sum_{m \neq l}^{d'} \frac{1}{\tau_{lm}^{c-}} \mathcal{C}_{lm}^2 \right) + \beta \sum_{l=1}^{d'} \Omega(\tau_l^{c+}) - \beta \sum_{l=1}^{d'} \sum_{m \neq l}^{d'} \Omega(\tau_{lm}^{c-})$	
LGP, CAN	$\ell_{\text{InfoNCE}}^\dagger = \sum_{c=1}^C \left(-\frac{1}{\tau_i^{c+}} \text{sim}(\mathbf{z}_i^c, \mathbf{z}_i^{c+}) + \frac{1}{\tau_{in*}^{c-}} \max_{n=1, \dots, N+1} \text{sim}(\mathbf{z}_i^c, \mathbf{z}_{in}^{c-}) \right) + \beta \Omega(\tau_i^{c+}) - \beta \Omega(\tau_{in*}^{c-})$	

The pipeline: further details



- The baseline uses one projection head and constant temperature, while our method has multiple projection heads and adaptive temperature.
- We use SimCLR for pre-training with ResNet-18 on STL-10. After pre-training, we choose 500 positive pairs and 500 negative pairs from the validation to compute the cosine similarity.
- In (a) and (b), similarity score (temperature scaled) is computed between the 128-dim features extracted from the projection head(s).
- In (c) and (d), cosine similarity score is computed between the 512-dim features extracted from the backbone after removing the project heads.



Results

