

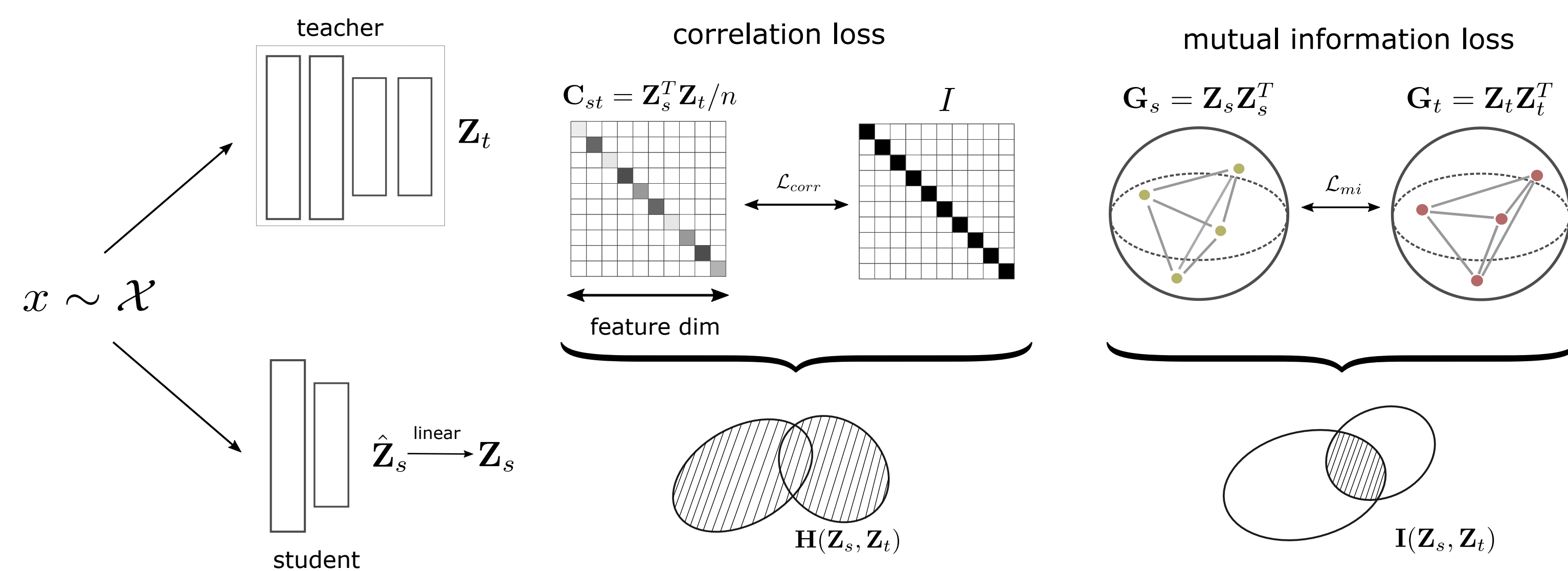


## Motivation

1. Most SoTA distillation methods are far **too computationally expensive** to adopt in practice, both in terms of their memory consumption, and FLOPs.
2. We introduce a **theoretical framework** for knowledge distillation that is rooted in information theory.
3. In doing so, we derive two complimentary losses that provide a **new SoTA on standard distillation benchmarks**.
4. We also consider two very different task, namely NLP, and binary network classification, thus demonstrating the **flexibility and scope** of our proposed losses.

## Method

Our proposed are derived from a set of **cheap** matrix-based estimators [1, 2] resembling Rényi's entropy in the single and multi-variate case.



**Correlation Loss** The correlation loss aims to match the information present in each **feature dimension** between the teacher and student representations. The parameter  $\alpha$  is related to Rényi's entropy order.

$$\mathcal{L}_{corr} = \log_2 \sum_{i=1}^d |v_i - 1|^{2\alpha} \quad (1)$$

**Mutual Information Loss** The mutual information loss provides an additional complimentary objective whereby we transfer the **intra-batch similarity** (i.e., the relationship between samples) from the teacher representations to the student representations. The second loss transfers to relationship between different data points within the batch.

$$\mathcal{L}_{mi} = \log_2 \|\mathbf{G}_s\|_F^2 - \log_2 \|\mathbf{G}_{st}\|_F^2 \quad (2)$$

The final loss is then just a weighted sum of these two. Experiments found that the model performance is **robust to the choice in loss weights**.

$$\mathcal{L}_{ITRD} = \mathcal{L}_{XE} + \beta_{corr} \mathcal{L}_{corr} + \beta_{mi} \mathcal{L}_{mi} \quad (3)$$

## Experiments & Results

CIFAR-100 test *accuracy* (%) of student networks trained with a number of distillation methods. The best results are highlighted in **bold**, while the second best results are underlined. ITRD achieves the best performance for 10 out of 13 of the architecture pairs, with a **6.8% and 24.4% relative improvement** over ReviewKD and WCoRD respectively.

| Teacher Student | W40-2 W16-2                             | W40-2 W40-1        | R56 R20            | R110 R20     | R110 R32           | R32x4 R8x4         | V13 V8             | V13 MN2            | R50 MN2            | R50 V8             | R32x4 SN1          | R32x4 SN2          | W40-2 SN1          |
|-----------------|---|--------------------|--------------------|--------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| Teacher Student | 75.61 73.26                             | 75.61 71.98        | 72.32 69.06        | 74.31 69.06  | 74.31 71.14        | 79.42 72.50        | 74.64 70.36        | 74.64 64.60        | 79.34 64.60        | 79.34 70.36        | 79.42 70.50        | 79.42 71.82        | 75.61 70.50        |
| KD              | 74.92                                   | 73.54              | 70.66              | 70.67        | 73.08              | 73.33              | 72.98              | 67.37              | 67.35              | 73.81              | 74.07              | 74.45              | 74.83              |
| CRD             | 75.64                                   | 74.38              | 71.63              | 71.56        | 73.75              | 75.46              | 74.29              | 69.94              | 69.54              | 74.58              | 75.12              | 76.05              | 76.27              |
| WCoRD           | <u>76.11</u>                            | <u>74.72</u>       | <b>71.92</b>       | <u>71.88</u> | <u>74.20</u>       | <u>76.15</u>       | 74.72              | 70.02              | 70.12              | 74.68              | 75.77              | 76.48              | 76.68              |
| ReviewKD        | <b>76.12</b>                            | <u>75.09</u>       | <u>71.89</u>       | -            | 73.89              | 75.63              | 74.84              | <u>70.37</u>       | 69.89              | -                  | <b>77.45</b>       | <b>77.78</b>       | <u>77.14</u>       |
| Ours            | $\mathcal{L}_{corr}$                    | 75.85 ±0.12        | 74.90 ±0.29        | 71.45 ±0.21  | 71.77 ±0.34        | 74.02 ±0.27        | 75.63 ±0.09        | 74.70 ±0.27        | 69.97 ±0.33        | <b>71.41</b> ±0.41 | <b>75.71</b> ±0.02 | 76.80 ±0.28        | 77.27 ±0.25        |
|                 | $\mathcal{L}_{corr} + \mathcal{L}_{mi}$ | <b>76.12</b> ±0.04 | <b>75.18</b> ±0.22 | 71.47 ±0.07  | <b>71.99</b> ±0.46 | <b>74.26</b> ±0.05 | <b>76.19</b> ±0.22 | <b>74.93</b> ±0.12 | <b>70.39</b> ±0.39 | <u>71.34</u> ±0.33 | <u>75.49</u> ±0.32 | <u>76.91</u> ±0.19 | <u>77.40</u> ±0.06 |

Experiments are also performed on **ImageNet**, where ITRD is comparable to state-of-the-art. The only two methods with a higher accuracy are 2× and 20× more computationally expensive respectively.

| Model          | EM         | F1          |             |
|----------------|------------|-------------|-------------|
| Teacher (BERT) | 81.5       | 88.6        |             |
| T6             | DistilBERT | 79.1        | 86.9        |
|                | TextBrewer | 80.8        | 88.1        |
|                | ITRD       | <b>81.5</b> | <b>88.5</b> |
| T3             | TextBrewer | 76.3        | 84.8        |
|                | ITRD       | <b>77.7</b> | <b>85.8</b> |

To show the wide applicability of our method, we consider distillation on a reading comprehension task. ITRD **outperforms both NLP-specific distillation methods TextBrewer and DistilBERT** in both the Exact Match (EM) metrics and in F1 score.

ITRD can be used to **reduce the gap between binary and full-precision (FP) networks**. Both CRD and ReviewKD degrade the BNN performance and, in contrast, ITRD improves upon the original ReCU by 1.3%, which is only 0.7% shy of the FP model.

| Network   | Method   | W/A   | Top-1 (%)   |
|-----------|--|-------|-------------|
| ResNet-18 | FP   | 32/32 | 94.8        |
|           | RBNN   | 1/1   | 92.2        |
|           | ReCU   | 1/1   | 92.8        |
|           | ReCU + CRD                                     | 1/1   | 92.1        |
|           | ReCU + ReviewKD                                | 1/1   | 92.6        |
|           | ReCU + $\mathcal{L}_{corr} + \mathcal{L}_{mi}$ | 1/1   | <b>94.1</b> |

ITRD losses can be implemented in **as few as 15 lines of code**.

We have also publicly releases the complete training and inference pipelines.

## References

- [1] Luis Gonzalo Sanchez Giraldo, Murali Rao, and Jose C. Principe. Measures of entropy from data using infinitely divisible Kernels. IEEE Transactions on Information Theory, 2015.
- [2] Paul L. Williams and Randall D. Beer. Nonnegative Decomposition of Multivariate Information. 2010.