# Functional Knowledge Transfer with Self-supervised Representation Learning

Prakash Chandra Chhipa[1,*], Muskaan Chopra[2,a], Gopal Mengi[2,b], Varun Gupta[2], Richa Upadhyay[1],
Meenakshi Subhash Chippa[1,c], Kanjar De[1], Rajkumar Saini[1], Seiichi Uchida[3] and Marcus Liwicki[1]

[1] *Luleå Tekniska Universitet, Luleå, Sweden*

{*prakash.chandra.chhipa, richa.upadhyay, kanjar.de, rajkumar.saini, marcus.liwicki*}*@ltu.se*

[c] meechi-2@student.ltu.se

[2] *CCET, Punjab University, Chandigarh, India*

{[a] *co19342,* [b] *co20320, varungupta*}*@ccet.ac.in*

[3]*Human Interface Laboratory, Kyushu University, Fukuoka, Japan*

*uchida@ait.kyushu-u.ac.jp*

[*]*Corresponding author - prakash.chandra.chhipa@ltu.se*

*Abstract*—This work investigates the unexplored usability of self-supervised representation learning in the direction of functional knowledge transfer. In this work, functional knowledge transfer is achieved by joint optimization of self-supervised learning pseudo task and supervised learning task, improving supervised learning task performance. Recent progress in self-supervised learning uses a large volume of data, which becomes a constraint for its applications on small-scale datasets. This work shares a simple yet effective joint training framework that reinforces human-supervised task learning by learning self-supervised representations just-in-time and vice versa. Experiments on three public datasets from different visual domains, Intel Image, CIFAR, and APTOS, reveal a consistent track of performance improvements on classification tasks during joint optimization. Qualitative analysis also supports the robustness of learnt representations. Source code and trained models are available on GitHub [1].
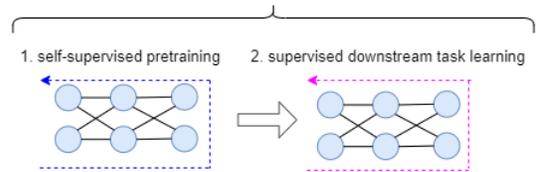
*Index Terms*—self-supervised learning, functional knowledge transfer, joint training, representation learning, computer vision
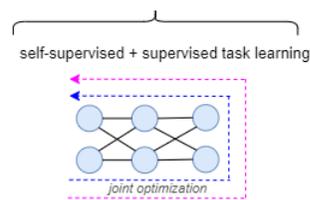
Fig. 1. Figure compares the proposed functional knowledge transfer approach in context of self-supervised learning and supervised task learning with conventional representational knowledge transfer approach where self-supervised pretraining and supervised task learning is performed in sequential manner.

## I. INTRODUCTION

The concept of functional knowledge transfer [1] has been explored for multi-task learning problems in computer vision [2]–[4] in the context of simultaneous training and joint optimization of multiple tasks. Typically functional knowledge transfer is employed for end-to-end joint training and optimization of multiple supervised learning tasks. Representational knowledge transfer, where pretraining and downstream task learning is done sequentially, has been thoroughly investigated and shown success in self-supervised learning. So far, functional knowledge transfer in self-supervised learning has not been studied, leaving a research gap.

This study uses functional knowledge transfer between self-supervised representation learning and other supervised downstream tasks. Figure 1 compares both knowledge transfer approaches. The proposed method jointly optimizes contrastive self-supervised learning with classification task learning on

ResNet-50 [5] backbone, explored on three public datasets of different visual domains, CIFAR10 [6], Intel Image [7], and Aptos [8]. The proposed approach enhances supervised task performance on all three datasets, supporting the hypothesis. Quantitative and qualitative comparisons are made between the proposed and conventional knowledge transfer approach. The following are the main contributions of this work:

1) Explored functional knowledge transfer with self-supervised representation learning towards making it applicable to the small-batch size and small-scale dataset.
2) Hypothesizes that self-supervised learning reinforces supervised task learning and vice versa.

With these contributions, proposed approach improves supervised task performance on all three datasets, supported by qualitative results and provide preliminary empirical support for the hypothesis.

---

[1]https://github.com/prakashchhipa/Functional_Knowledge_Transfer_SSL

## II. RELATED WORK

Joint embedding architecture and method based self-supervised learning has shown significant advances in label-free representation learning paradigm. It is based on learning similarity in transformed views of input images and the way it learns robust features by avoiding collapsed representation it is divided into several categories, e.g., i) Contrastive Methods (SimCLR [9], MoCo [10]), ii) Distillation (BYOL [11], Sim-Siam [12]), iii) Clustering (SwAV [13]), and (iv) Information Maximization (Barlow Twins [14], VICReg [15]). All these methods have explored the representational knowledge transfer approach, where pretraining is performed, and learned parameters are transferred as knowledge to enable downstream tasks. However, functional knowledge transfer and simultaneous training are unexplored. Although some work has been carried out to exploit the label details in self-supervised methods [16], especially contrastive learning.

On the other side, multi-task learning [2]–[4], [17] has explored functional knowledge transfer by simultaneous training procedures is their natural requirement and has shown progress toward improved performance and computational efficiency. Self-supervised learning approaches for functional knowledge transfer are unexplored. It could make self-supervised algorithms computationally efficient and adaptable to small datasets by integrating with other learning tasks.

## III. METHOD

The proposed method enables a specific type of inductive transfer, called functional knowledge transfer [1] on self-supervised representation learning approach by incorporating simultaneous training with downstream task learning. Specifically, the proposed method employs the contrastive learning method [9] for self-supervised representation learning and classification as downstream tasks on multiple datasets, CIFAR10 [6], Aptos [8], and Intel Image [7]. The following section describes the method in detail.

Data $D : (X, Y)$ is set of input sample pair of $(x, y)$ where $x \in \mathbb{R}^d$, is the input image data of $d$ dimensions and $y$ is corresponding human-annotation from annotation space $\mathcal{C}$. The data is defined as $D : \{(x_1, y_1), ...(x_n, y_n)\} \subseteq \mathbb{R}^d \times \mathcal{C}$.

### A. Contrastive Self-supervised Learning

To define the joint embedding architecture and method based self-supervised learning objective, followed in contrastive learning (SimCLR [9]), a set of $K$ non-leanrable transformations $\mathcal{T} : \{t_k\}_{k \in K}$ is defined, which are image processing based augmentations, provides transformed views of input image $(x', x'')$, to retain the invariant feature learning. Further, learnable function $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ parameterized by learnable parameters $\Theta_f$ which is Convolutional Neural Network (CNN) backbone and another learnable function $g : \mathbb{R}^m \rightarrow \mathbb{R}^{\bar{m}}$ parameterized by learnable parameters $\Theta_g$ which is projector network is defined. With that, Noise contrastive estimation [18] based self-supervised contrastive

learning objective, NT-Xent (Normalized Temperature Scaled Cross Entropy) loss is defined in Eq. 1.

$$\mathcal{L}_{SSL} = \sum_{(x', x'') \in \mathcal{T}(X)} - \log \frac{e^{\mathcal{A}}}{\sum_{k=1}^{2|X|} 1_{[k \neq x']} e^{\mathcal{B}}} \quad (1)$$

$$\mathcal{A} = (sim(g(\Theta_g; f(\Theta_f; x')), g(\Theta_g; f(\Theta_f; x''))))/\tau \quad (2)$$

$$\mathcal{B} = (sim(g(\Theta_g; f(\Theta_f; x')), g(\Theta_g; f(\Theta_f; x^k))))/\tau \quad (3)$$

where, $\mathcal{A}$ defines similarity for positive pairs, $\mathcal{B}$ constitute similarity for negative pairs with denominator part of Eq. 1, and $sim$ is cosine similarity and $\tau$ is temperature scale parameter, and $\mathcal{L}_{SSL}$ is contrastive loss.

### B. Supervised Task Learning

Supervised learning objective for mentioned downstream task of classification can be mentioned in terms of cross entropy loss $\mathcal{L}_{CE}$, defined in Eq. 4.

$$\mathcal{L}_{CE} = -\frac{1}{|D|} \sum_{(x,y) \in D} \sum_{c \in \mathcal{C}} y_c \log(f(\Theta; x_c)) \quad (4)$$

### C. Representational Knowledge Transfer

Representational Knowledge Transfer is extensively explored in self-supervised learning, not only in contrastive learning but also in other self-supervised paradigms, i.e., distillation [11], [12] and information maximization [14], [15]. This type of knowledge transfer comprises two stages;

1) First stage is self-supervised pretraining of CNN backbone without requiring labels which learns invariant representations of underlying visual concepts by similarity learning, described in Eq. 1
2) Second stage is downstream supervised tasks learning in which learnt representations from stage one is used by initializing the learning parameters of CNN encoder, and supervised training is performed accordance to the task, e.g., classification, described in Eq. 4

The first part of Figure 1 symbolically depicts the process.

### D. Functional Knowledge Transfer

Functional Knowledge Transfer in the context of self-supervised learning is defined by jointly optimizing the self-supervised learning objective with supervised task learning objective. $\mathcal{L}_{FKT}$ loss described in Eq. 5 is single stage process where parameters learning is simultaneous and influenced by both loss objectives in just-in-time manner. $\lambda$ is parameter for balancing losses, however kept 1 in all experiments.

$$\mathcal{L}_{FKT} = \mathcal{L}_{CE} + \lambda \, \mathcal{L}_{SSL} \quad (5)$$

The second part of Figure 1 symbolically demonstrate the process.

*Analytical Reasoning*: Functional Knowledge Transfer in context of self-supervised learning with supervised task learning is based on reinforced effects of tasks to each other. More concretely, it is shown in Figure 3 and defined as:

- Invariant Features - Self-supervised learning objective shares invariant generalized descriminative features,
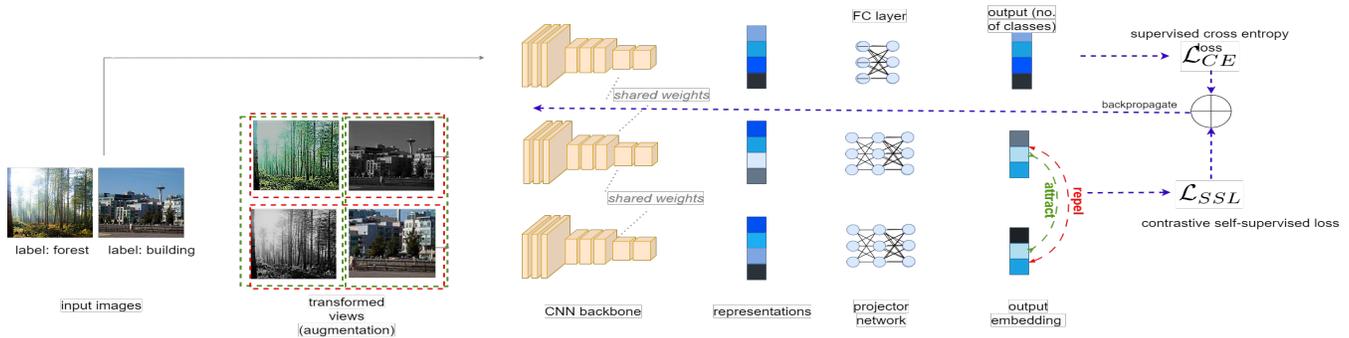
Fig. 2. Illustrates the Functional Knowledge Transfer where contrastive loss and cross entropy loss is computed on self-supervised and supervised tasks respectively and jointly backpropagated, which enables simultaneous training.



Fig. 3. Demonstrate the bi-directional constructive reinforcements for self-supervised learning and supervised task learning which enables self-supervision on relatively smaller batch size and small-scale datasets and improves classification performance.

which reinforces the task specific feature learning for given human annotation

- Robust Semantics - Supervised task learning shares robust semantic information (e.g., categorization, clusters of similar concepts) of underlying visual concepts of image backed by human knowledge, which reinforces similarity learning of semantically similar visual concepts

This bi-directional constructive reinforcements improves learning of both the tasks, which can enable to learn contrastive learning on relatively smaller batch sizes and smaller datasets and improved performance for supervised downstream task, shown in Figure 3.

## IV. DATASETS

This study uses public datasets of natural geographic scenes, atomic objects, and medical images to investigate functional knowledge transfer on self-supervised representation learning in diverse visual concepts. The Table I summarizes the three datasets used in this work.

TABLE I
DATASETS

| Dataest | Image type | No. of images | | No. of |
|---|---|---|---|---|
| | | Train | Test | classes |
| CIFAR-10 [6] | singular objects | 50000 | 10000 | 10 |
| Intel Image [7] | natural scenes | 14034 | 3000 | 6 |
| APTOS 2019 [8] | retinal images | 3263 | 399 | 5 |

## V. EXPERIMENTAL DETAILS

To evaluate the applicability of contrastive self-supervised learning method in functional knowledge transfer approach, detailed experimentation was performed on three public datasets, CIFAR10 [6], Intel Image [7], and Aptos [8] from diverse visual domains. Functional knowledge transfer is employed by joint training of self-supervised (simCLR [9]) and supervised task learning (classification) as mentioned in the Section III. A comparative study is performed by benchmarking the proposed approach to conventional approach of representational knowledge transfer, where the model is pretrained and then trained for the downstream task.

Methodological investigations are preferred; hence, common hyperparameters are configured for all three datasets with both transfer knowledge approaches. To emphasize a less compute-intensive approach, single GPU implementation is preferred with ResNet-50 [5] backbone and batch size of 256 for contrastive learning, which is much smaller than the original work. Due to this very reason, contrastive pretraining on CIFAR is perfomed with batch size 256, which was not available elsewhere. Pretraining, downstream task, and joint training are configured for 100 epochs. Self-supervised pretraining in both approaches uses LARS optimizer with learning rate 0.001 temperature scale 0.5 and employs standard augmentations suggested in the original work simCLR [9]. Supervised learning classification tasks in both approaches use SGD optimizer with a learning rate 0.025. All the experiments are repeated three times and the mean value of the performance metric is reported with standard deviation.

## VI. RESULTS AND DISCUSSIONS

Table II describes the multi-class classification performance of the proposed approach by comparing it with the conventional approach for all three datasets. A consistent improvement, up to $1.40\%$, is observed in accuracy for all three datasets for the proposed functional knowledge transfer approach. It is worth noting that all the results show negligible standard deviation across several trial. The proposed approach has improved performance over previous work on APTOS and intel image datasets, also supported by qualitative analysis. Important observations are briefly described as follows -

**Functional Knowledge Transfer improves performance**: Results comparisons in Table II clearly show the inspiring trend that functional transfer has improved the downstream

| Dataset | Method | Accuracy | Precision | Recall |
|---|---|---|---|---|
| CIFAR10 | Representational Transfer $ | 92.20±0.11 | 92.18±0.10 | 92.21±0.10 |
| | Functional Transfer | **93.60±0.10** | **93.62±0.13** | **93.59±0.11** |
| Intel Image | Representational Transfer | 93.18±0.15 | 93.15±0.18 | 93.17±0.20 |
| | Functional Transfer | **93.70±0.13** | **93.33±0.11** | **93.31±0.11** |
| Aptos 2019 | Representational Transfer | 83.10±0.10 | 83.05±0.09 | **83.05±0.12** |
| | Functional Transfer | **83.32±0.11** | **83.14±0.10** | 83.04±0.10 |

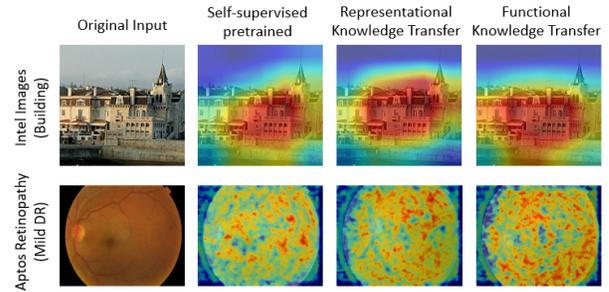| Dataset | Method | Accuracy | Precision | Recall |
|---|---|---|---|---|
| Intel Image | Representational Transfer | 78.28±0.14 | 78.16±0.15 | 78.19±0.18 |
| | Functional Transfer | 78.52±0.13 | 78.38±0.14 | 78.43±0.12 |



Fig. 4. Pretrained model, representational knowledge transfer, and functional transfer approaches are compared for class activation maps (CAM). First instance is from building category from intel image dataset, and second instance is mild DR category from APTOS dataset. CAM not produced for CIFAR10 dataset due to very small size of input.
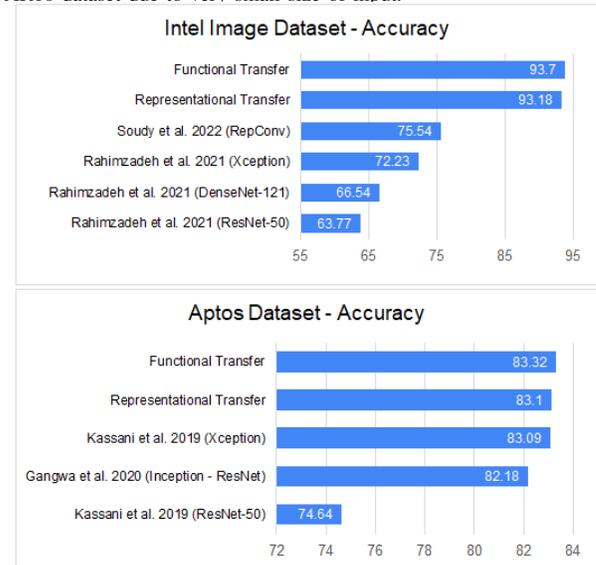


Fig. 5. Comparison with previous works: Intel Image (top), Aptos retinopathy fundus (bottom).

task performance regardless of the dataset over the conventional approach. It also outperforms previous works, using the same ResNet-50 architecture and beyond, as shown in Figure 5 for APTOS and Intel Image datasets.

**Functional Knowledge Transfer enables efficient self-supervision**: Enabling self-supervised learning on small-scale datasets and smaller batch size is another significant outcome for the functional knowledge transfer approach. It supports the hypothesis mentioned in Figure 3 where both tasks reinforced the efficiency to each other. However, more investigation is required in efficiently fusing self-supervised and supervised learning task loss objectives for even further improved performance.

**Functional Knowledge Transfer demonstrates computational efficiency**: Representational knowledge transfer requires 100 epochs of pretraining followed by 100 epochs of downstream supervised task learning. In contrast, the functional knowledge transfer approach performs better by joint training for 100 epochs. Effectively, functional knowledge transfer requires roughly half of computations or at-least saves downstream task computation costs. It is also essential to evaluate the functional knowledge transfer for domain adaptation and other transfer learning scenarios in future work because self-supervised learning in representation knowledge transfer intends to do the transfer learning.

**Qualitative Robustness**: Quantitative results and performance are also supported by qualitative analysis shown through class activation maps in Figure 4 for two datasets, Intel Image and Aptos, where attention regions are displayed.

It clearly shows the ability to attend to the region of interest to capture the essence of visual concepts in the images. When compared to the pretrained and representational knowledge transfer approaches, functional knowledge transfer demonstrated very competitive and even more focused attention region.

**Ablation**: Ablation study is performed on ResNet-18 backbone on Intel Image dataset (Table III), which shows marginal improvement, which gives motivation to investigate further in this direction.

## VII. CONCLUSION

Functional knowledge transfer is explored on contrastive self-supervised learning with classification tasks where exciting performance improvement is depicted across multiple public datasets. It has shown preliminary empirical support for enabling contrastive self-supervised learning on small batches and small-scale datasets by reinforcing the task during joint training. This study strongly encourages further investigation of functional knowledge transfer using different self-supervised learning paradigms and supervised learning tasks.

## REFERENCES

[1] R. Vilalta, C. Carrier, P. Brazdil, C. M. Soares *et al.*, "Inductive transfer," 2017.

[2] R. Caruana, "A dozen tricks with multitask learning," in *Neural networks: tricks of the trade*. Springer, 2002, pp. 165–191.

[3] A. Maurer, M. Pontil, and B. Romera-Paredes, "The benefit of multitask representation learning," *Journal of Machine Learning Research*, vol. 17, no. 81, pp. 1–32, 2016.

[4] M. Crawshaw, "Multi-task learning with deep neural networks: A survey," *arXiv preprint arXiv:2009.09796*, 2020.

[5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[6] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

[7] I. corp, "image classification challenge," 2019.

[8] S. D. Karthik and, Maggie, "Aptos 2019 blindness detection," 2019.

[9] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.

[10] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9729–9738.

[11] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar *et al.*, "Bootstrap your own latent-a new approach to self-supervised learning," *Advances in neural information processing systems*, vol. 33, pp. 21 271–21 284, 2020.

[12] X. Chen and K. He, "Exploring simple siamese representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 15 750–15 758.

[13] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," *Advances in neural information processing systems*, vol. 33, pp. 9912–9924, 2020.

[14] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, "Barlow twins: Self-supervised learning via redundancy reduction," in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 310–12 320.

[15] A. Bardes, J. Ponce, and Y. LeCun, "Vicreg: Variance-invariance-covariance regularization for self-supervised learning," *arXiv preprint arXiv:2105.04906*, 2021.

[16] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," *Advances in neural information processing systems*, vol. 33, pp. 18 661–18 673, 2020.

[17] R. Upadhyay, P. C. Chhipa, R. Phlypo, R. Saini, and M. Liwicki, "Multi-task meta learning: learn how to adapt to unseen tasks," *arXiv preprint arXiv:2210.06989*, 2022.

[18] M. Gutmann and A. Hyvärinen, "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 297–304.