

Advanced Machine Learning

Paper: *Learning without forgetting* [13]

Romain Mormont

Montefiore Institute, ULiège

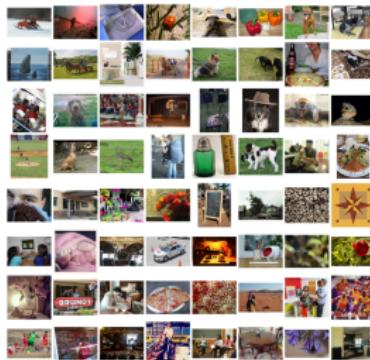
April 2, 2019

Image classification and transfer learning

Deep learning works well with images...
as long as we have enough of them !

What if we have a small dataset ?

⇒ **transfer learning** from a larger database (e.g. ImageNet) !



Transfer learning

"Transfer learning aims to produce an effective model for a target task with limited or no labeled training data by leveraging and exploiting knowledge from a different, but related source domain to predict the truth label for an unseen target instance."

[3]

We consider homogeneous transfer learning:

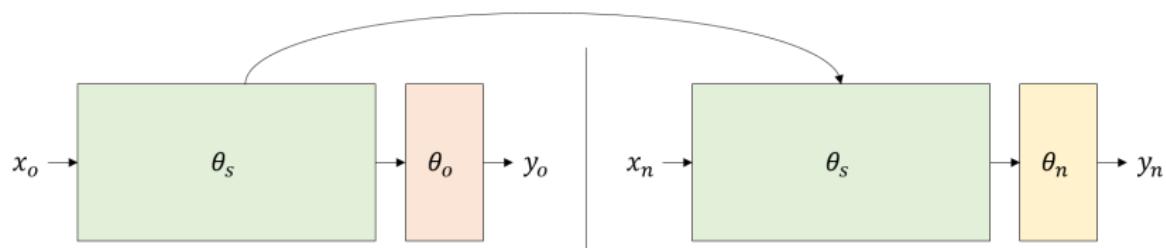
- a **source** task $(\mathcal{X}_o, \mathcal{Y}_o, p_o(x, y))$ (**old** task)
- a **target** task $(\mathcal{X}_n, \mathcal{Y}_n, p_n(x, y))$ (**new** task)
- homogeneity: $\mathcal{X}_o = \mathcal{X}_n$ (images) but $p(x_o) \neq p(x_n)$
- different tasks: $\mathcal{Y}_o \neq \mathcal{Y}_n$
- usually: $N_n \ll N_o$

Transfer learning

We consider the following network components:

- θ_s : shared parameters
- θ_o : source task-specific parameters
- θ_n : target task-specific parameters

θ_s and θ_o are **pre-trained on the source task** and θ_n is **initialized randomly**.



How to transfer ?

There exists several transfer techniques:

- **feature extraction:** θ_s and θ_o are unchanged. Features are extracted from one or more layers to train θ_n or another classifier.
- **fine-tuning:** both θ_s and θ_n are optimized on the target task (θ_o fixed).
- **partial fine-tuning:** similar to fine-tuning but early layers of the shared network are frozen. Only a part of θ_s is optimized along with θ_n .
- **joint/multitask training:** θ_s , θ_n and θ_o are all optimized with both tasks at the same time (e.g. by interleaving samples).
- ...

Feature extraction I

How to transfer ?

θ_s and θ_o are unchanged. Features are extracted from one or more layers to train θ_n or another classifier.

The idea of extracting features from pre-trained networks arrived not so long after the deep learning revolution:

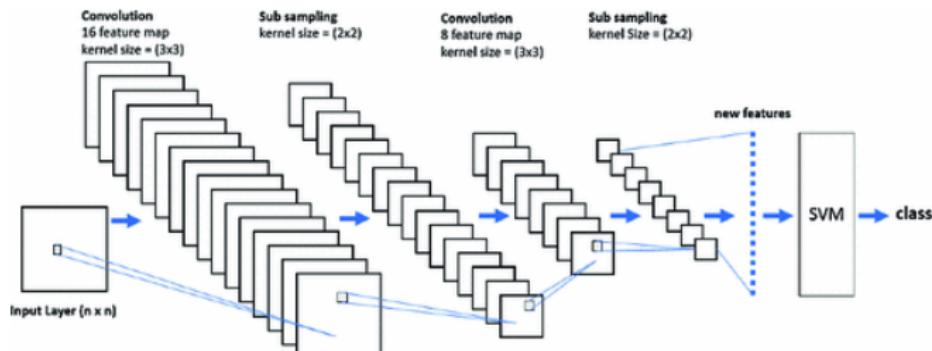
- Decaf [5] (2013.10): feature extractor based on AlexNet
- Overfeat [15] (2013.12): another feature extractor based on AlexNet

Nowadays, all modern deep learning frameworks provide pre-trained weights for **well-known architectures** (e.g. ResNet, DenseNet, GoogLeNet,...) which can then be used as feature extractors.

Feature extraction II

One usually extracts **one or more vector of features** per input image.

The extracted vectors can then be used to train classifier with **traditional machine learning** methods (e.g. SVM, ANN, random forests,...). Simple linear models are often enough to obtain competitive performances.



Feature extraction III

Advantages

- no need for training the network \Rightarrow fast to apply
- benefit from complex features learnt by the transferred network
- deep features often yield better performances than handcrafted features

Drawbacks

- features are not task-specific \Rightarrow **under-performing**

Conclusion

Feature extraction provides a **good baseline** for image classification and transfer learning.

Fine-tuning I

How to transfer ?

Both θ_s and θ_n are optimized on the target task (θ_o fixed).

Fine-tuning has been applied with great success in many applications and fields [6, 7].

All parameters are trained to **minimize a loss on a new task**, usually with a **small learning rate**. Fine-tuning is a way of making the features more discriminative for the new task. Using a small learning rate is a way of preserving the representations learnt from the source task.

Fine-tuning II

Advantages

- accelerate convergence at training
- learned features are more task-specific
- often **outperform features extraction** or learning from random weights

Drawbacks

- subject to **catastrophic forgetting**

Conclusion

One of the most efficient approaches when tackling a new classification task.

Catastrophic forgetting

"[Catastrophic forgetting] refers to the catastrophic loss of previously learned responses, whenever an attempt is made to train the network with a single new (additional) response. Affected networks include, for example, backpropagation learning networks."

[2]

Catastrophic forgetting (CF):

- a neural network has no built-in mechanism to retain previously learnt information
- when fine-tuning, CF is actually a consequence of overfitting the new task [8]
- mechanisms aimed at handling this problem would actually act as regularization during training and would **reduce overfitting**

Less/no catastrophic forgetting means better transfer !

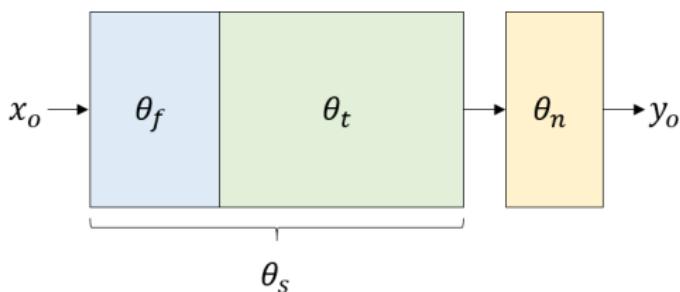
Partial fine-tuning I

How to transfer ?

Similar to fine-tuning but early layers of the shared network are frozen.
Only a part of θ_s is optimized along with θ_n .

Let $\theta_f \subseteq \theta_s$ be the set of frozen parameters and $\theta_t = \theta_s \setminus \theta_f$ be the set of trainable parameters. Partial fine-tuning is a **compromise between features extraction and fine-tuning**, as it falls back to:

- fine-tuning when $\theta_f = \emptyset$
- features extraction when $\theta_f = \theta_s$



Partial fine-tuning II

Similarly to fine-tuning, the parameters θ_t and θ_n are trained to **minimize a loss on the new task.**

Partial fine-tuning offers a mechanism for mitigating **catastrophic forgetting and overfitting** by fixing a part of the network.

However, it doesn't solve the problem especially when few parameters are frozen.

Partial fine-tuning III

Advantages

- mechanism for dealing with catastrophic forgetting
- can accelerate training, as less parameters to optimize

Drawbacks

- the optimal extent of θ_f is task-specific
- catastrophic forgetting is not avoided

Conclusion

Partial fine-tuning offers a mechanism for dealing with catastrophic forgetting but is not the final solution

Joint training I

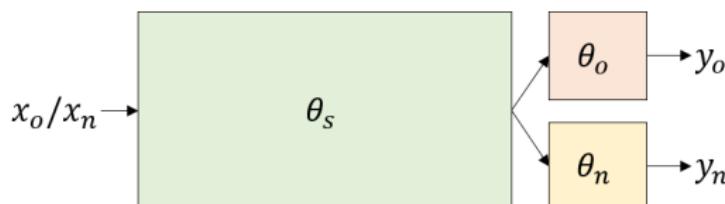
How to transfer ?

θ_s , θ_n and θ_o are all optimized with both tasks at the same time.

Also known as **multi-task training** [1]. Both tasks are provided at the same time (instead of sequentially).

Regarding training, several questions arise:

- interleaving samples ? interleaving tasks ?
- how to combine task-specific losses ?
- what about imbalance when $N_o \gg N_n$?



Joint training II

Advantages

- catastrophic forgetting is not an issue, or is greatly mitigated

Drawbacks

- more data \Rightarrow longer training time
- require access to the (potentially large) source dataset

Conclusion

Joint training deals with catastrophic forgetting at the expense of longer training time and is not possible if the source dataset is unavailable.

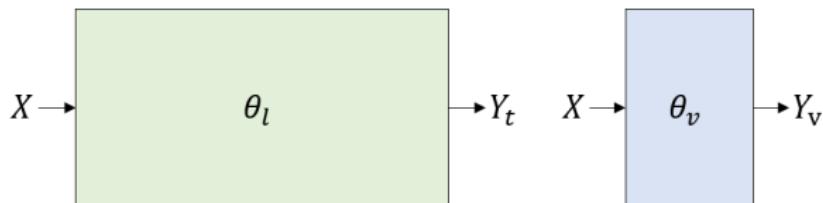
Summary

	Feature Extraction	Fine Tuning	Joint Training
new task performance	X medium	good	best
original task performance	good	X bad	good
training efficiency	fast	fast	X slow
storage requirement	medium	medium	X large
requires previous task data	no	no	X yes

Note: knowledge distillation

There exists single task transfer techniques that consists in transferring knowledge learnt by a large network θ_l into a significantly smaller network θ_v with $|\theta_l| \gg |\theta_v|$.

One of such techniques is called **distillation** [9] and consists in optimizing the small network **to replicate the probabilities** produced by the larger network for a same input sample using a distillation loss.



Learning without Forgetting

Zhizhong Li, Derek Hoiem, *Member, IEEE*

Abstract—When building a unified vision system or gradually adding new capabilities to a system, the usual assumption is that training data for all tasks is always available. However, as the number of tasks grows, storing and retraining on such data becomes infeasible. A new problem arises where we add new capabilities to a Convolutional Neural Network (CNN), but the training data for its existing capabilities are unavailable. We propose our Learning without Forgetting method, which uses only new task data to train the network while preserving the original capabilities. Our method performs favorably compared to commonly used feature extraction and fine-tuning adaption techniques and performs similarly to multitask learning that uses original task data we assume unavailable. A more surprising observation is that Learning without Forgetting may be able to replace fine-tuning with similar old and new task datasets for improved new task performance.

Index Terms—Convolutional Neural Networks, Transfer Learning, Multi-task Learning, Deep Learning, Visual Recognition

1 INTRODUCTION

MANY practical vision applications require learning new visual capabilities while maintaining performance on existing ones. For example, a robot may be delivered to someone's house with a set of default object recognition capabilities, but new site-specific object models need to be added. Or for construction safety, a system can identify whether a worker is wearing a safety vest or hard hat, but a superintendent may wish to add the ability to detect improper footware. Ideally, the new tasks could be learned while sharing parameters from old ones, without

network could be duplicated and fine-tuned for each new task to create a set of specialized networks.

It is also possible to use a variation of fine-tuning where part of θ_n – the convolutional layers – are frozen to prevent overfitting, and only top fully connected layers are fine-tuned. This can be seen as a compromise between fine-tuning and feature extraction. In this work we call this method **Fine-tuning FC** where FC stands for fully connected.

Joint Training (e.g. [7]) All parameters $\theta_1, \theta_2, \theta_n$ are

Learning without forgetting

Timeline:

2016.09 - First publication on arXiv

2017.11 - Published in *IEEE Transactions on Pattern Analysis and Machine Intelligence, Volume 40*

Summary:

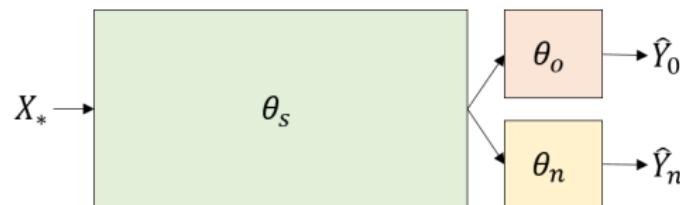
- transfer learning for image classification with deep convolutional network
- new method for training a pre-trained network on a new task
- **goal**: avoid "*forgetting*" what was learnt on the source task
- **idea**: a loss function designed to preserve the network response for the old task during training on the new task

Method I

Inputs:

- shared parameters θ_s and θ_o pre-trained on the source task
- target task data X_n and Y_n
- new layer(s) with parameters θ_n attached to the last layer of the shared network

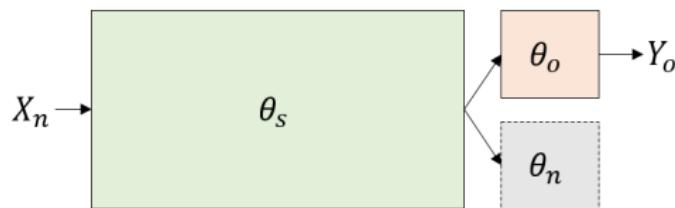
No need for source task data !



Method II

Initialization:

- extraction of "old" network responses Y_o (i.e. old task classes probabilities/softmaxes) to the new data X_n
- θ_n parameters random initialization



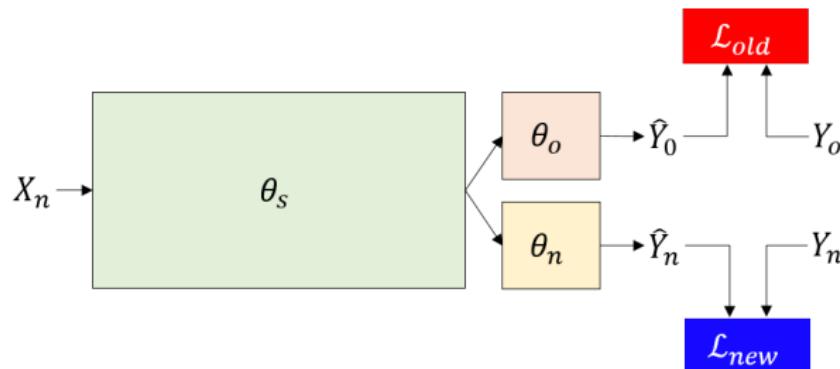
Method III

Training: jointly optimize θ_s , θ_o and θ_n with stochastic gradient descent:

$$\hat{\theta}_s^*, \hat{\theta}_o^*, \hat{\theta}_n^* = \underset{\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n}{\operatorname{argmin}} \left(\lambda_0 \mathcal{L}_{old} (Y_0, \hat{Y}_0) + \mathcal{L}_{new} (Y_n, \hat{Y}_n) + \mathcal{R} (\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n) \right)$$

where:

- $\lambda_0 \in \mathbb{R}$ is a loss balance weight
- \mathcal{L}_{old} is the old task loss
- \mathcal{L}_{new} is the new task loss
- \mathcal{R} is a regularization term (weight decay)

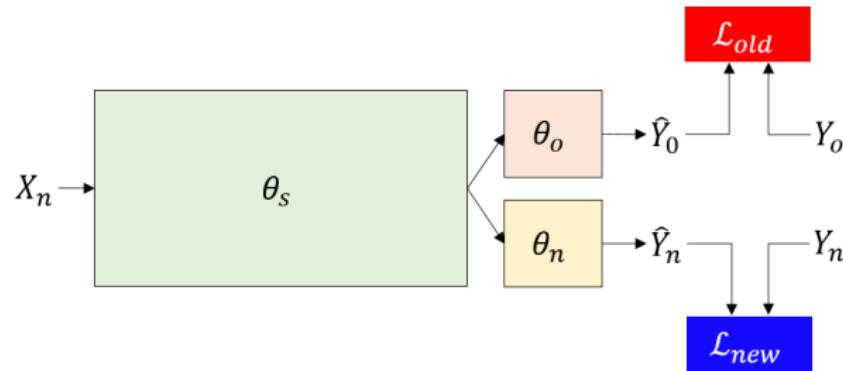


\mathcal{L}_{new} (new task loss)

Goal: should encourage predicted class probabilities to be consistent with ground truth. An obvious pick for obtaining this behavior is the **multinomial logistic loss** (averaged over a training batch):

$$\mathcal{L}_{new}(\mathbf{y}_n, \hat{\mathbf{y}}_n) = -\mathbf{y}_n \cdot \log \hat{\mathbf{y}}_n$$

where \mathbf{y}_n is the one-hot ground-truth class vector and $\hat{\mathbf{y}}_n$ is the softmax output of the network.



\mathcal{L}_{old} (old task loss)

Goal: output probabilities \hat{Y}_o should be close to the recorded outputs probabilities Y_0 . They use the distillation loss of [9]:

$$\mathcal{L}_{old}(\mathbf{y}_o, \hat{\mathbf{y}}_o) = - \sum_{i=1}^C y_0'^{(i)} \log \hat{y}_0'^{(i)}$$

where C is the number of classes and $y_0'^{(i)}$ and $\hat{y}_0'^{(i)}$ are a modified version of the softmax function:

$$y_0'^{(i)} = \frac{\left(y_o^{(i)}\right)^{1/T}}{\sum_j \left(y_o^{(j)}\right)^{1/T}}$$

$$\hat{y}_0'^{(i)} = \frac{\left(\hat{y}_o^{(i)}\right)^{1/T}}{\sum_j \left(\hat{y}_o^{(j)}\right)^{1/T}}$$

where T is the "*temperature*" hyper-parameter of the distillation loss.

Algorithm

LEARNING WITHOUT FORGETTING:

Start with:

θ_s : shared parameters

θ_o : task specific parameters for each old task

X_n, Y_n : training data and ground truth on the new task

Initialize:

$Y_o \leftarrow \text{CNN}(X_n, \theta_s, \theta_o)$ // compute output of old tasks for new data
 $\theta_n \leftarrow \text{RANDINIT}(|\theta_n|)$ // randomly initialize new parameters

Train:

Define $\hat{Y}_o \equiv \text{CNN}(X_n, \hat{\theta}_s, \hat{\theta}_o)$ // old task output

Define $\hat{Y}_n \equiv \text{CNN}(X_n, \hat{\theta}_s, \hat{\theta}_n)$ // new task output

$\theta_s^*, \theta_o^*, \theta_n^* \leftarrow \underset{\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n}{\operatorname{argmin}} \left(\lambda_o \mathcal{L}_{old}(Y_o, \hat{Y}_o) + \mathcal{L}_{new}(Y_n, \hat{Y}_n) + \mathcal{R}(\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n) \right)$

Similar method: Less-Forgetting Learning [10]

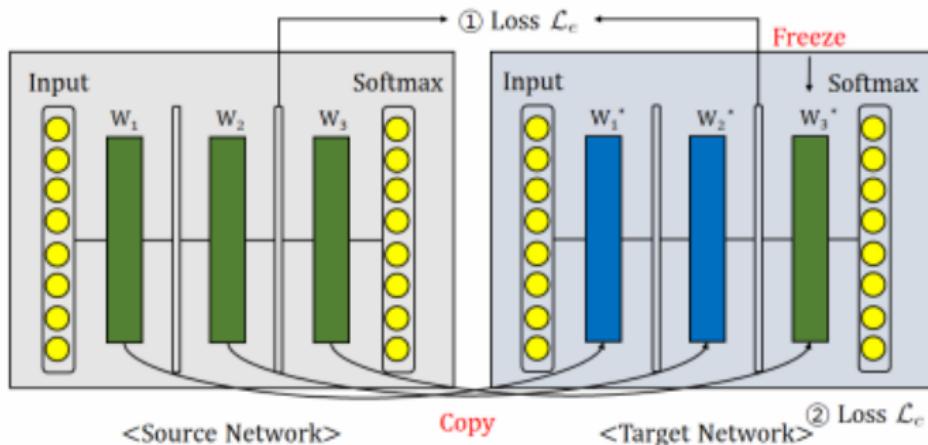


Fig. 2. Conceptual diagram for describing a less-forgetting method. Our learning method uses the trained weights of the source network as the initial weights of the target network and minimizes two loss functions simultaneously.

Experiments

Their experiments aim at evaluate whether *Learning without Forgetting* is an **effective method for learning a new task while preserving performances on old tasks.**

They compare their method in various settings to:

- features extraction
- fine-tuning
- fine-tuning FC (i.e. partial fine-tuning)
- joint-training (**upper bound**)
- Less Forgetting Learning (LFL)

To study the algorithms, they use several datasets.

- As source: ImageNet [4] and Places365 [18]
- As target: Pascal VOC2012 everingham2015pascal, Caltech-UCSD Bird [17], MIT Indoor Scenes [14] and MNIST [12]

Experiments - compared methods

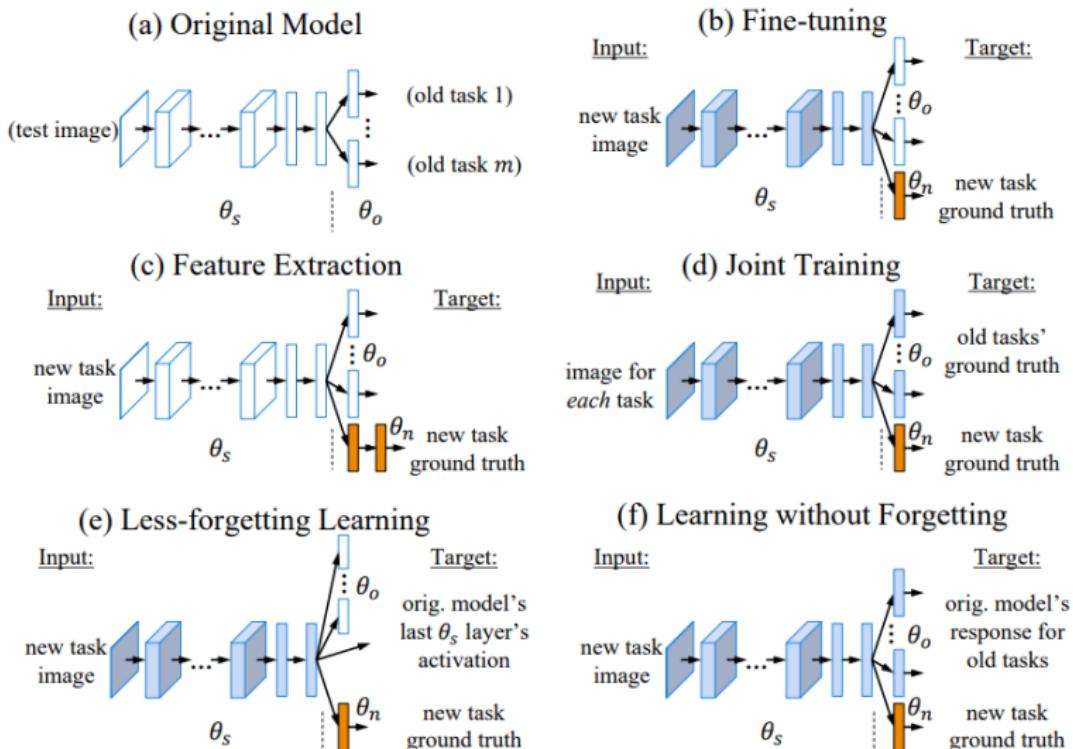
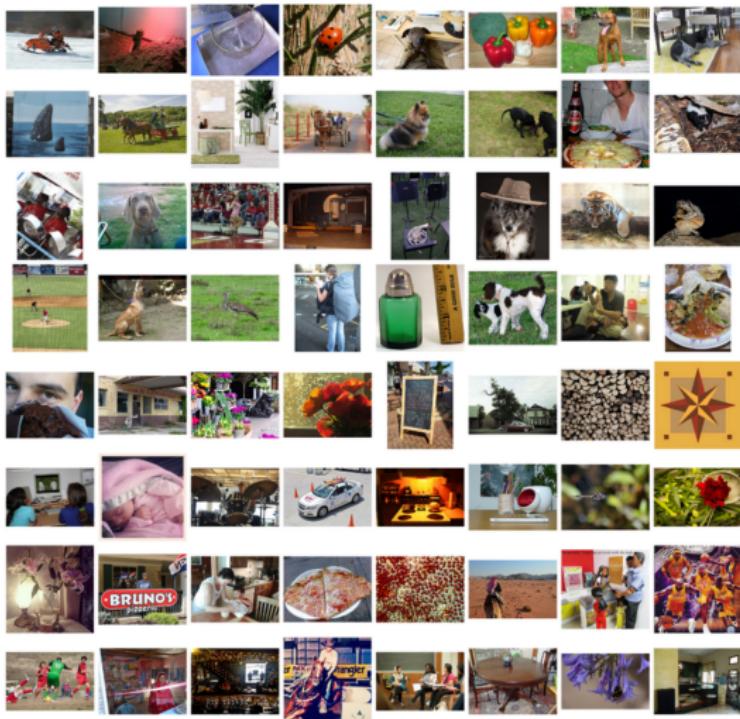


Fig. 2. Illustration for our method (f) and methods we compare to (b-e). Images and labels used in training are shown. Data for different tasks are used in alternation in joint training.

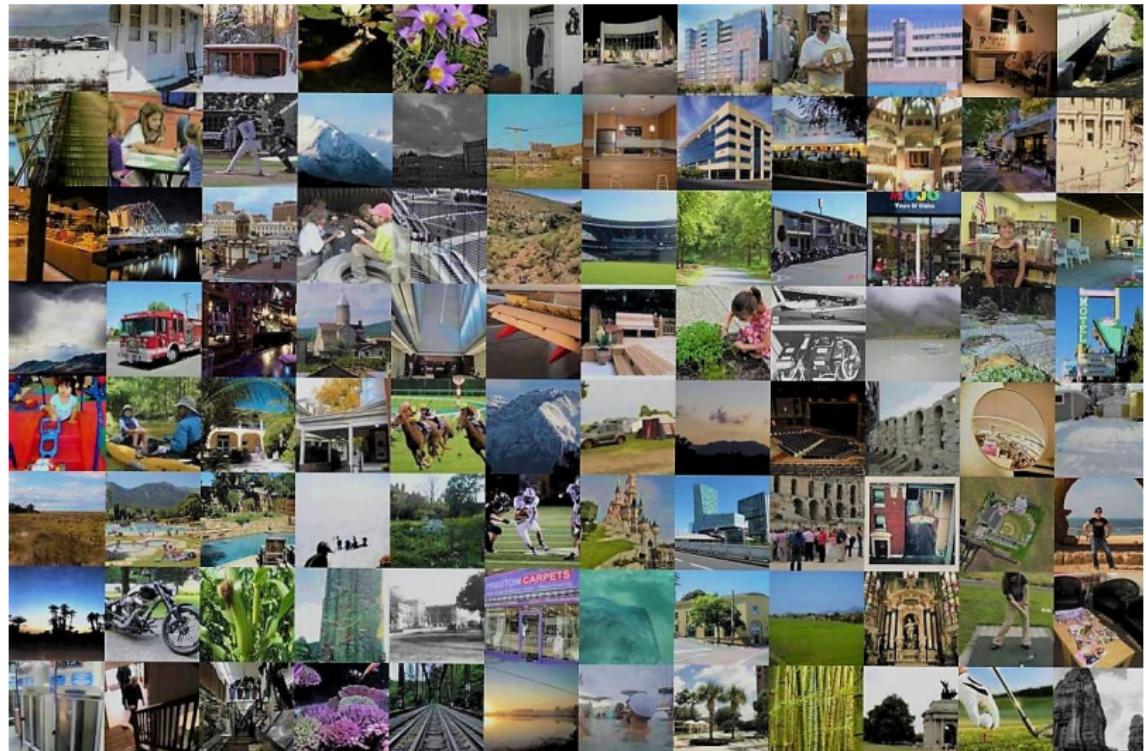
Datasets - ImageNet

> 1 million images, 1000 classes



Datasets - Places365

1.6 million images, 365 classes



Datasets - Pascal VOC 2012

5717 images, similar to ImageNet



Aeroplanes



Bicycles



Birds



Boats



Bottles



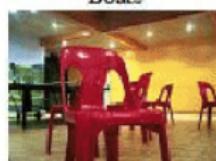
Buses



Cars



Cats



Chairs



Cows



Dining tables



Dogs



Horses



Motorbikes



People



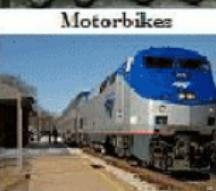
Potted plants



Sheep



Sofas



Trains



TV/Monitors

Datasets - CUB

5994 images, bird recognition



Datasets - Scenes

5360 images, indoor scenes



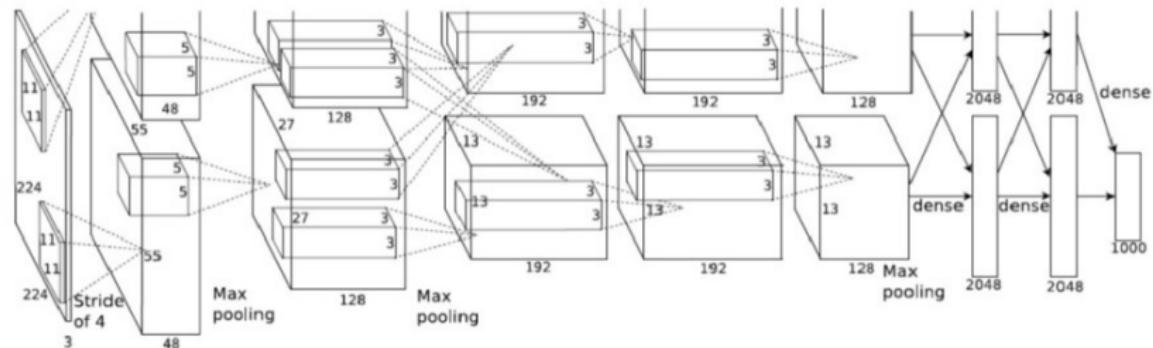
Datasets - MNIST

60000 images, digit recognition



Experiments - networks

AlexNet [11], 62M parameters:



VGG 16 [16], 138M parameters:



<https://goo.gl/U1oCWy>

<https://goo.gl/4UHmwv>

Experiments - "single task scenario" I

Settings: learning one new task among different task pairs and different methods.

	ImageNet→VOC		ImageNet→CUB		ImageNet→Scenes		ImageNet→MNIST	
	old	new	old	new	old	new	old	new
LwF (ours)	56.2	76.1	54.7	57.7	55.9	64.5	49.8	99.3
Fine-tuning	-0.9	-0.3	-3.8	-0.7	-2.0	-0.8	-2.8	0.0
LFL	0.0	-0.4	-1.9	-2.6	-0.3	-0.9	-2.9	-0.6
Fine-tune fc	0.5	-0.7	0.2	-3.9	0.6	-2.1	7.0	-0.2
Feat. Extraction	0.8	-0.5	2.3	-5.2	1.2	-3.3	7.3	-0.8
Joint Training	0.7	-0.2	0.6	-1.1	0.5	-0.6	7.2	-0.0
	Places365→VOC		Places365→CUB		Places365→Scenes		Places365→MNIST	
	old	new	old	new	old	new	old	new
LwF (ours)	50.6	70.2	47.9	34.8	50.9	75.2	38.3	99.2
Fine-tuning	-2.2	0.1	-4.6	1.0	-2.1	-1.7	-0.9	0.1
LFL	0.2	-0.7	0.7	-1.7	-0.2	-0.5	-0.4	-0.1
Fine-tune fc	0.5	-1.3	1.8	-4.9	0.3	-1.1	13.0	-0.2
Feat. Extraction	1.1	-1.4	3.8	-12.3	0.8	-1.7	13.3	-1.1
Joint Training	0.7	-0.0	2.3	1.5	0.3	-0.3	13.4	-0.1

"Performance for the single new task scenario. For all tables, the difference of methods' performance with LwF (our method) is reported to facilitate comparison. Mean Average Precision is reported for VOC and accuracy for all others."

Experiments - "single task scenario" II

Settings: learning one new task among different task pairs and different methods.

	ImageNet→VOC		ImageNet→CUB		ImageNet→Scenes		ImageNet→MNIST	
	old	new	old	new	old	new	old	new
LwF (ours)	56.2	76.1	54.7	57.7	55.9	64.5	49.8	99.3
Fine-tuning	-0.9	-0.3	-3.8	-0.7	-2.0	-0.8	-2.8	0.0
LFL	0.0	-0.4	-1.9	-2.6	-0.3	-0.9	-2.9	-0.6
Fine-tune fc	0.5	-0.7	0.2	-3.9	0.6	-2.1	7.0	-0.2
Feat. Extraction	0.8	-0.5	2.3	-5.2	1.2	-3.3	7.3	-0.8
Joint Training	0.7	-0.2	0.6	-1.1	0.5	-0.6	7.2	-0.0
	Places365→VOC		Places365→CUB		Places365→Scenes		Places365→MNIST	
	old	new	old	new	old	new	old	new
LwF (ours)	50.6	70.2	47.9	34.8	50.9	75.2	38.3	99.2
Fine-tuning	-2.2	0.1	-4.6	1.0	-2.1	-1.7	-0.9	0.1
LFL	0.2	-0.7	0.7	-1.7	-0.2	-0.5	-0.4	-0.1
Fine-tune fc	0.5	-1.3	1.8	-4.9	0.3	-1.1	13.0	-0.2
Feat. Extraction	1.1	-1.4	3.8	-12.3	0.8	-1.7	13.3	-1.1
Joint Training	0.7	-0.0	2.3	1.5	0.3	-0.3	13.4	-0.1

"On the **new task**, our method consistently outperforms LFL, fine-tuning FC and feature extraction, while outperforming fine-tuning on most task pairs"

"Dissimilar new tasks degrade old task performance more."

Experiments - "single task scenario" III

Settings: learning one new task among different task pairs and different methods.

	ImageNet→VOC		ImageNet→CUB		ImageNet→Scenes		ImageNet→MNIST	
	old	new	old	new	old	new	old	new
LwF (ours)	56.2	76.1	54.7	57.7	55.9	64.5	49.8	99.3
Fine-tuning	-0.9	-0.3	-3.8	-0.7	-2.0	-0.8	-2.8	0.0
LFL	0.0	-0.4	-1.9	-2.6	-0.3	-0.9	-2.9	-0.6
Fine-tune fc	0.5	-0.7	0.2	-3.9	0.6	-2.1	7.0	-0.2
Feat. Extraction	0.8	-0.5	2.3	-5.2	1.2	-3.3	7.3	-0.8
Joint Training	0.7	-0.2	0.6	-1.1	0.5	-0.6	7.2	-0.0
	Places365→VOC		Places365→CUB		Places365→Scenes		Places365→MNIST	
	old	new	old	new	old	new	old	new
LwF (ours)	50.6	70.2	47.9	34.8	50.9	75.2	38.3	99.2
Fine-tuning	-2.2	0.1	-4.6	1.0	-2.1	-1.7	-0.9	0.1
LFL	0.2	-0.7	0.7	-1.7	-0.2	-0.5	-0.4	-0.1
Fine-tune fc	0.5	-1.3	1.8	-4.9	0.3	-1.1	13.0	-0.2
Feat. Extraction	1.1	-1.4	3.8	-12.3	0.8	-1.7	13.3	-1.1
Joint Training	0.7	-0.0	2.3	1.5	0.3	-0.3	13.4	-0.1

"On the **old task**, our method performs better than fine-tuning but often underperforms feature extraction, fine-tuning FC, and occasionally LFL."

"Our method usually performs similarly to **joint training** with AlexNet"

Experiments - "single task scenario" IV

Settings: learning one new task among different task pairs and different methods.

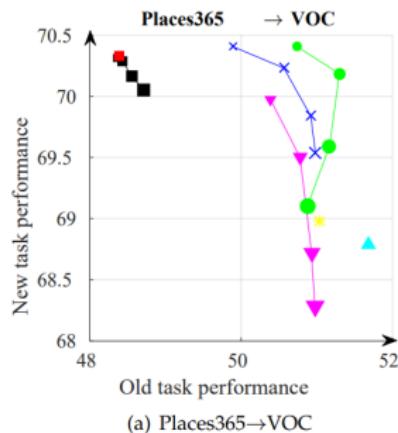
(c) Using VGG structure

	ImageNet→CUB		ImageNet→Scenes	
	old	new	old	new
LwF (ours)	60.6	72.5	66.8	74.9
Fine-tuning	-9.9	0.6	-4.1	-0.3
LFL	0.3	-2.8	-0.0	-2.1
Fine-tune fc	3.2	-6.7	1.4	-2.4
Feat. Extraction	8.2	-8.6	1.9	-5.1
Joint Training	8.0	2.5	4.1	1.5

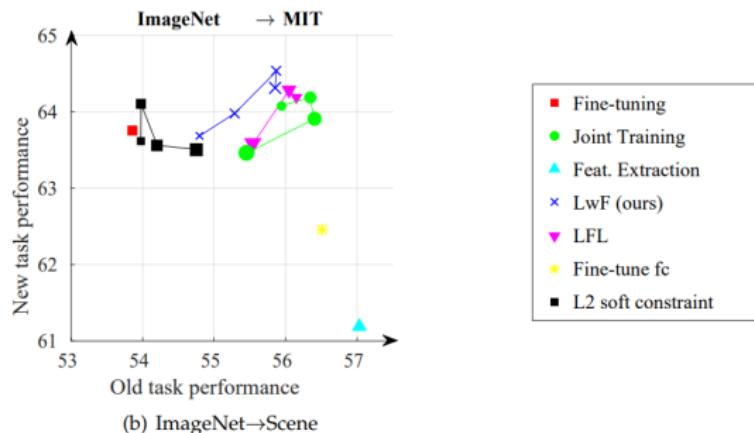
"[...] except that joint training outperforms consistently for VGG, and LwF performs worse than before on the old task. [...] joint training may have a larger benefit on **networks with more representational power.**"

Experiments - "single task scenario" V

Settings: learning one new task among different task pairs and different methods.



(a) Places365→VOC



(b) ImageNet→Scene

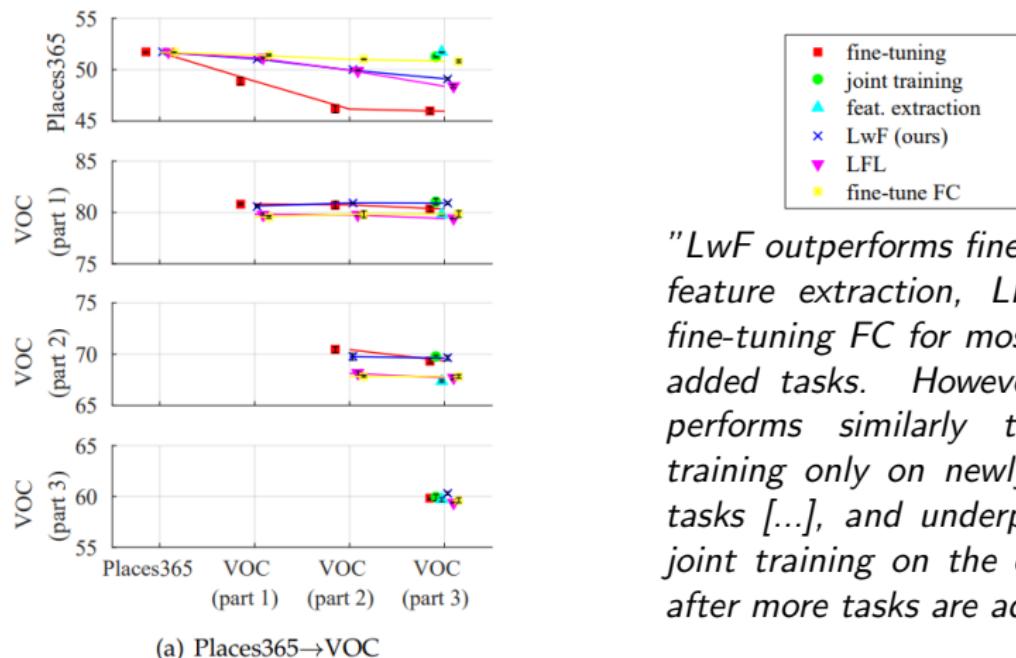
- Fine-tuning
- Joint Training
- ▲ Feat. Extraction
- ✖ LwF (ours)
- ▼ LFL
- ★ Fine-tune fc
- L2 soft constraint

"Visualization of both new and old task performance for compared methods, some with different weights of losses. [...] Larger symbols signifies larger λ_0 , i.e. heavier weight towards response-preserving loss."

"[...] if λ_0 is adjusted, LwF can perform better than LFL and fine-tuning FC on the new task for the same old task performance."

Experiments - "multiple new tasks scenario" |

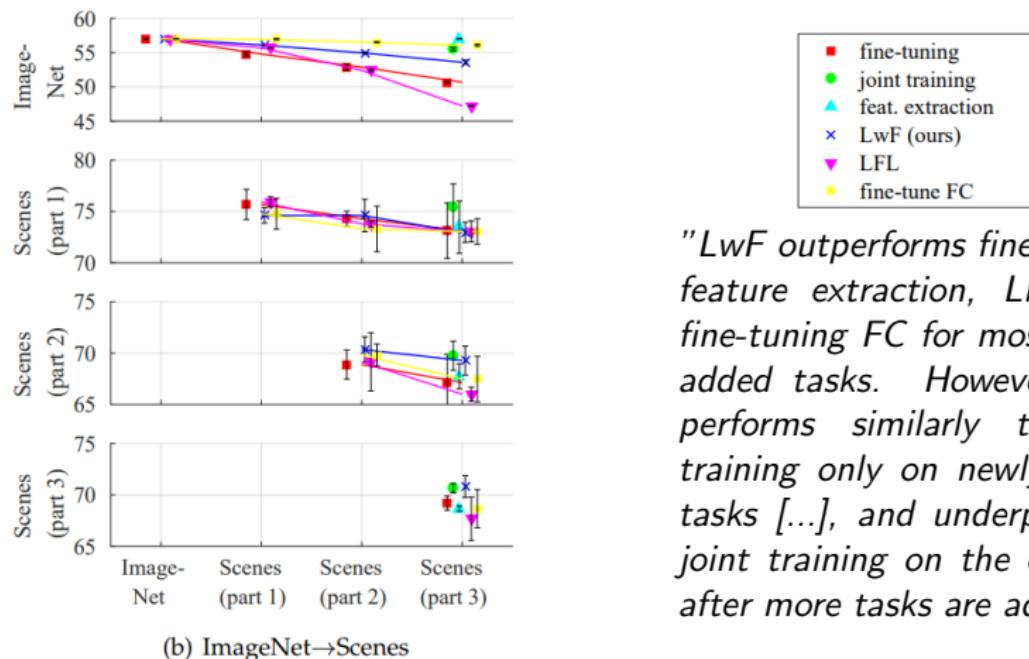
Settings: we compare different methods when cumulatively adding new tasks to the system.



"*LwF outperforms fine-tuning, feature extraction, LFL, and fine-tuning FC for most newly added tasks. However, LwF performs similarly to joint training only on newly added tasks [...], and underperforms joint training on the old task after more tasks are added.*"

Experiments - "multiple new tasks scenario" II

Settings: we compare different methods when cumulatively adding new tasks to the system.



"LwF outperforms fine-tuning, feature extraction, LFL, and fine-tuning FC for most newly added tasks. However, LwF performs similarly to joint training only on newly added tasks [...], and underperforms joint training on the old task after more tasks are added."

Limitations

Limitations:

- tasks must be enumerated: cannot deal with domains changing continuously on a spectrum
- not designed to support streaming data (new dataset should be entirely present before training)
- exhibit some *forgetting* (especially with particularly different tasks)
- ⇒ performance depends on the tasks " *relatedness*"
- with larger networks, performance gap between joint training and LwF becomes significant

Conclusion

Conclusion:

- A **new method for mitigating catastrophic forgetting** when training a pre-trained model on a new task
- **Good alternative to joint training (?)** which does not require access to the source data !
- **Good alternative to fine-tuning (?)**

Wrapping up - my opinion

Good & bads:

- ✓ simple and interesting method
- ✓ well-motivated, good review of related methods
- ✓ several varied dataset and several runs experiments (though only 3 runs)
- ✓ convincing single task experiment, numerous small experiments
- ✗ multiple tasks experiment unconvincing
- ✗ what about testing more networks with a larger capacity ?
- ✗ bold conclusion and title (still actually forgetting !)

I believe this method **should not replace joint training or fine-tuning** as implicitly suggested in the conclusion because the experiments are not complete enough to safely generalize and the method has issues with dissimilar tasks. A **method worth trying** however when working on such transfer problem !

References I

-  Rich Caruana. "Multitask learning". In: *Machine learning* 28.1 (1997), pp. 41–75.
-  *Catastrophic forgetting.*
<http://standoutpublishing.com/g/catastrophic-forgetting.html>. Accessed: 2019-03-19.
-  Oscar Day and Taghi M Khoshgoftaar. "A survey on heterogeneous transfer learning". In: *Journal of Big Data* 4.1 (2017), p. 29.
-  Jia Deng et al. "Imagenet: A large-scale hierarchical image database". In: *2009 IEEE conference on computer vision and pattern recognition*. ieee. 2009, pp. 248–255.
-  Jeff Donahue et al. "Decaf: A deep convolutional activation feature for generic visual recognition". In: *International conference on machine learning*. 2014, pp. 647–655.

References II

-  Andre Esteva et al. "Dermatologist-level classification of skin cancer with deep neural networks". In: *Nature* 542.7639 (2017), p. 115.
-  Varun Gulshan et al. "Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs". In: *Jama* 316.22 (2016), pp. 2402–2410.
-  Steven Gutstein and Ethan Stump. "Reduction of catastrophic forgetting with transfer learning and ternary output codes". In: *2015 International Joint Conference on Neural Networks (IJCNN)*. IEEE. 2015, pp. 1–8.
-  Geoffrey Hinton, Oriol Vinyals and Jeff Dean. "Distilling the knowledge in a neural network". In: *arXiv preprint arXiv:1503.02531* (2015).
-  Heechul Jung et al. "Less-forgetting learning in deep neural networks". In: *arXiv preprint arXiv:1607.00122* (2016).

References III

-  Alex Krizhevsky, Ilya Sutskever and Geoffrey E Hinton. "Imagenet classification with deep convolutional neural networks". In: *Advances in neural information processing systems*. 2012, pp. 1097–1105.
-  Yann LeCun et al. "Gradient-based learning applied to document recognition". In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.
-  Zhizhong Li and Derek Hoiem. "Learning without forgetting". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40.12 (2018), pp. 2935–2947.
-  Ariadna Quattoni and Antonio Torralba. "Recognizing indoor scenes". In: *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE. 2009, pp. 413–420.
-  Pierre Sermanet et al. "Overfeat: Integrated recognition, localization and detection using convolutional networks". In: *arXiv preprint arXiv:1312.6229* (2013).

References IV

-  Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition". In: *arXiv preprint arXiv:1409.1556* (2014).
-  Catherine Wah et al. "The caltech-ucsd birds-200-2011 dataset". In: (2011).
-  Bolei Zhou et al. "Places: An image database for deep scene understanding". In: *arXiv preprint arXiv:1610.02055* (2016).