

Day 4

Speaker: KanPeek Alahari , India

Title: Incremental Visual Learning

## Incremental Learning ?

- Continual learning
- Lifelong learning
- Sequential learning
- Never-ending Learning

} *Synonyms*

KA: Incremental Learning

3

## An Incremental Learning Scenario

- Growing up in India



*Overtime*

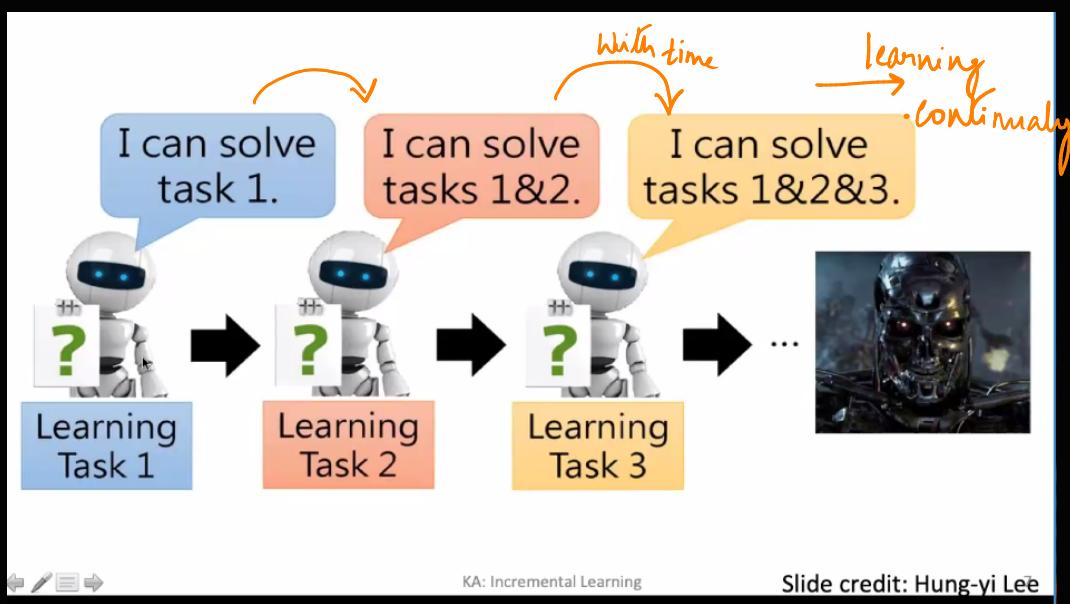
## An Incremental Learning Scenario

- And then during travels



KA: Incremental Learning

5



## Standard Machine Learning

**TRAIN – VALIDATION – TEST**

All sampled from the same distribution

-> benchmarks and academic datasets 😊 ✓ Able to do

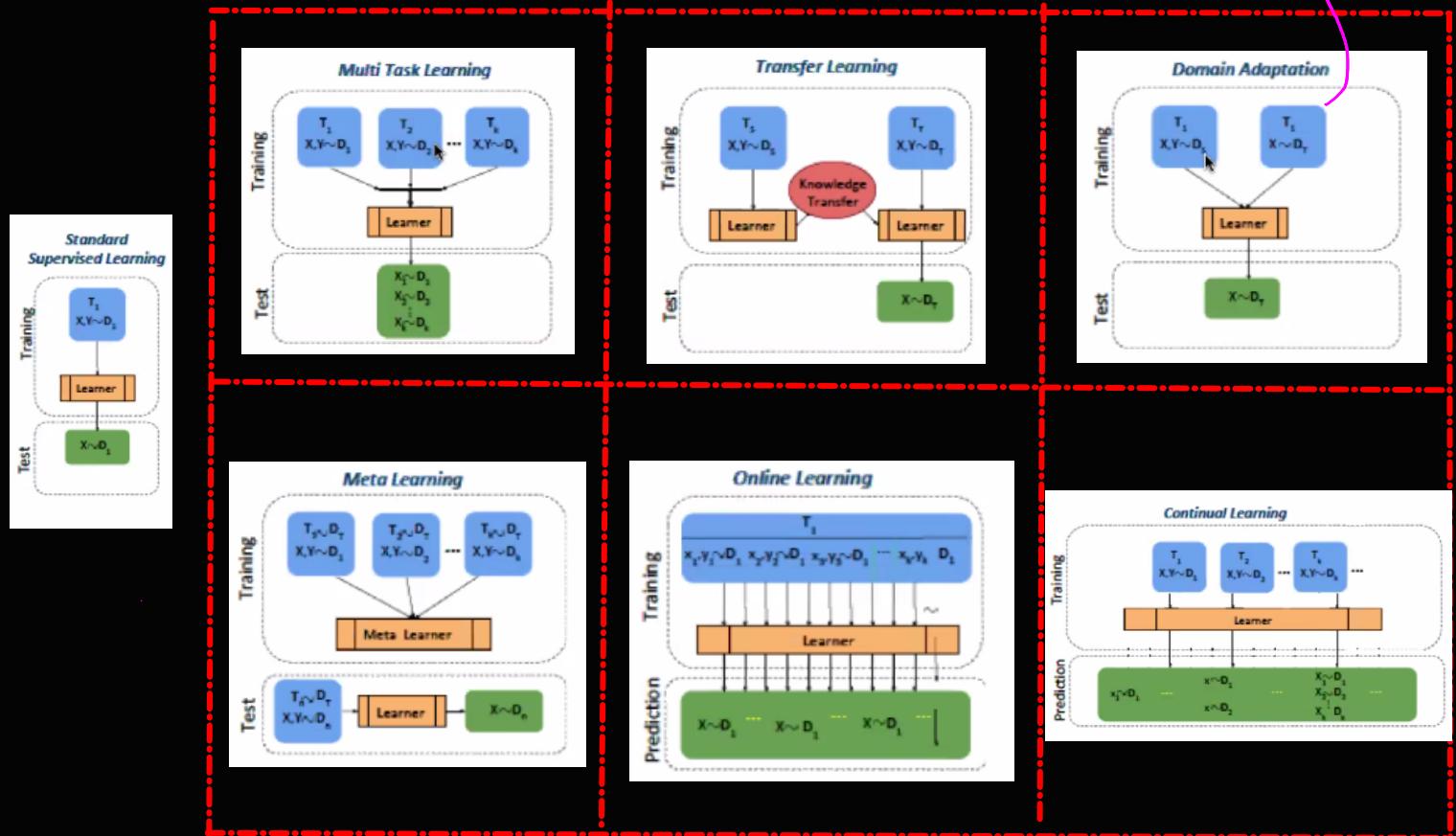
-> real-world systems ↗ Under Research

-> embodied learning ↗ Not yet.

Don't have all labelled images.

KA: Incremental Learning

Slide credit: T. Tuytelaars



# Incremental Learning Setup



- Task-incremental learning (*when  $T_1, T_2$  are different tasks*)
- Class-incremental learning
- Domain-incremental learning



KA: Incremental Learning

Slide credit: T. Tuytelaars

## Incremental Learning

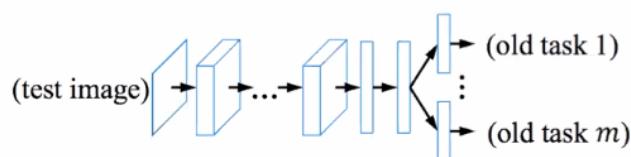
- A classical problem in machine learning, e.g.,  
[Carpenter et al. '92, Cauwenberghs and Poggio '00, Polikar et al. '01,  
Schlimmer and Fisher '86, Thrun '96]
- Some methods
  - Zero-shot learning, e.g., [Lampert et al. '13]  
**No training step for unseen classes**
  - Continuously update the training set, e.g., [Chen et al. '13]  
**Keep data and retrain**
  - Use a fixed data representation, e.g., [Mensink et al. '13]  
**Simplify the learning problem**

KA: Incremental Learning

12

## Brute Force Solution (non-incremental)

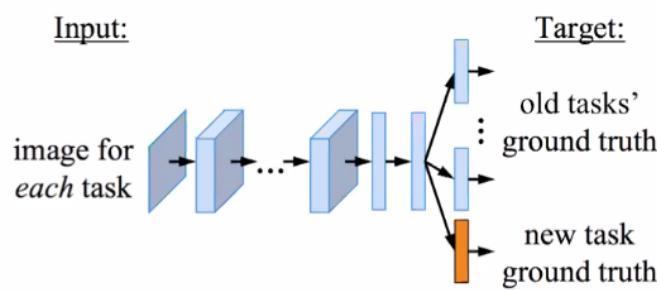
Original model



Joint training

"golden" baseline

- random initialize + train
- fine-tune
- unchanged



KA: Incremental Learning

Figures from [Li and Hoiem 2016]

# Brute Force Solution (non-incremental)

Retrain full model with both old and new data

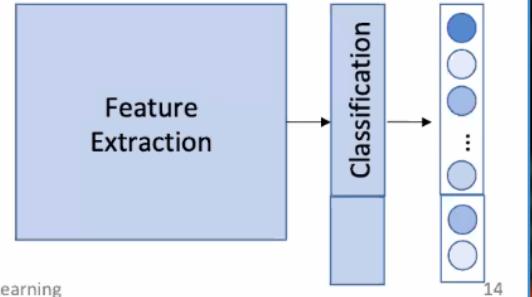
- Computationally expensive

eg.  
Video data.  
Healthcare  
data  
Storage Issue

- Needs access to old data
  - Storage capacity limitations
  - Privacy issues
  - Scalability issues

Slide credit: T. Tuytelaars

KA: Incremental Learning



14

## Why not brute force ?

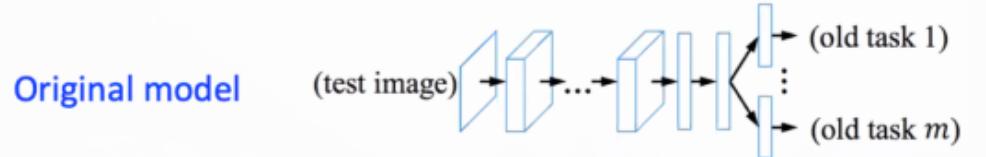
Sensitive data {  
large dataset {  
Data in the previous model  
maybe inaccessible.

- No access to all the data
- Can not store all the data
- Access to only a previously learned model, e.g., trained by others

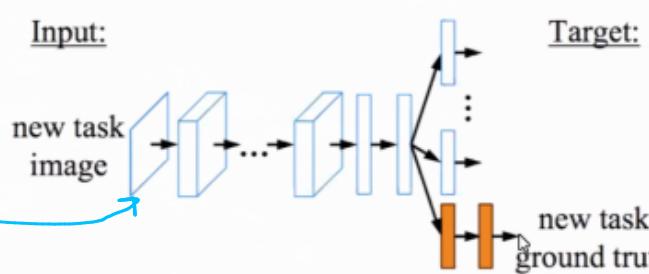
KA: Incremental Learning

15

## Naïve Solution 1



Feature extraction



random initialize + train

fine-tune

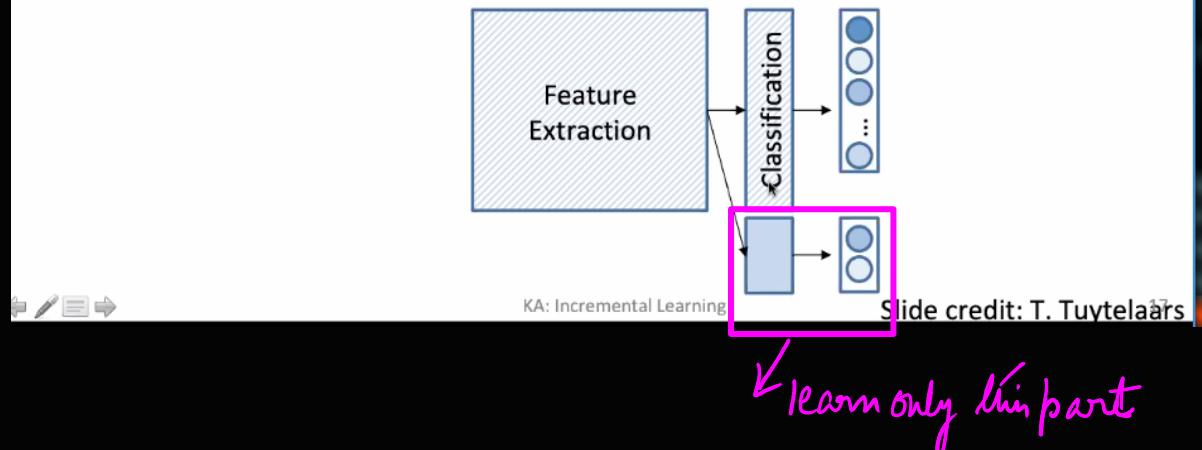
unchanged

KA: Incremental Learning

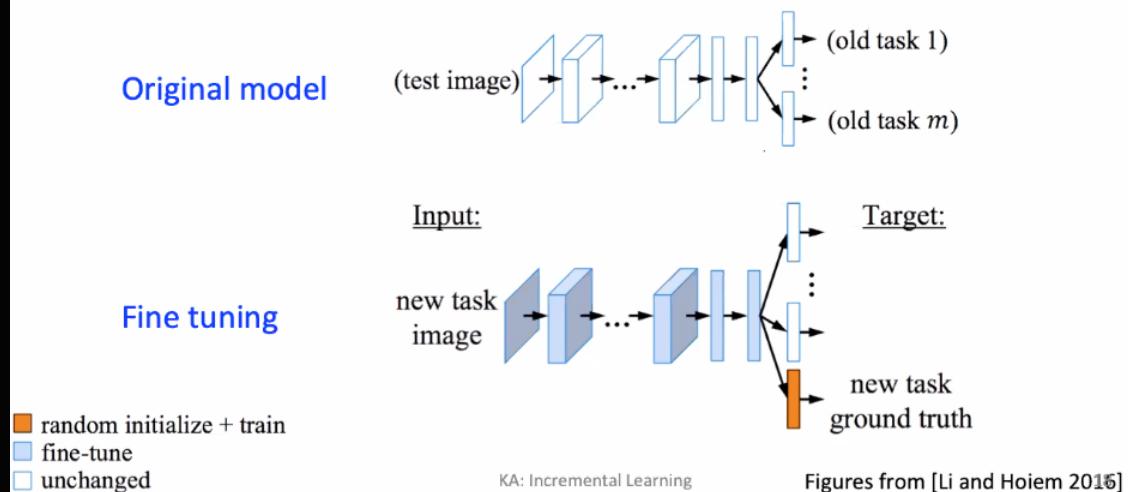
Figures from [Li and Hoiem 2014]

## Naïve Solution 1

- Finetune only last layer using new data only
- Suboptimal Results

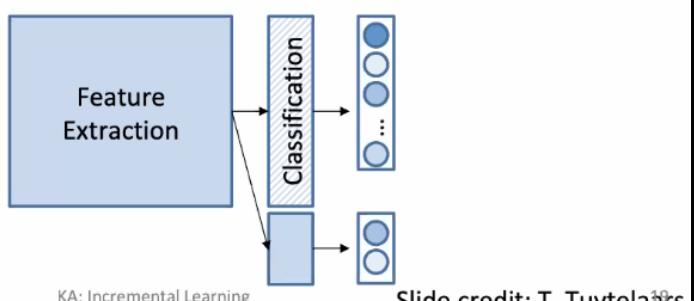


## Naïve Solution 2



## Naïve Solution 2

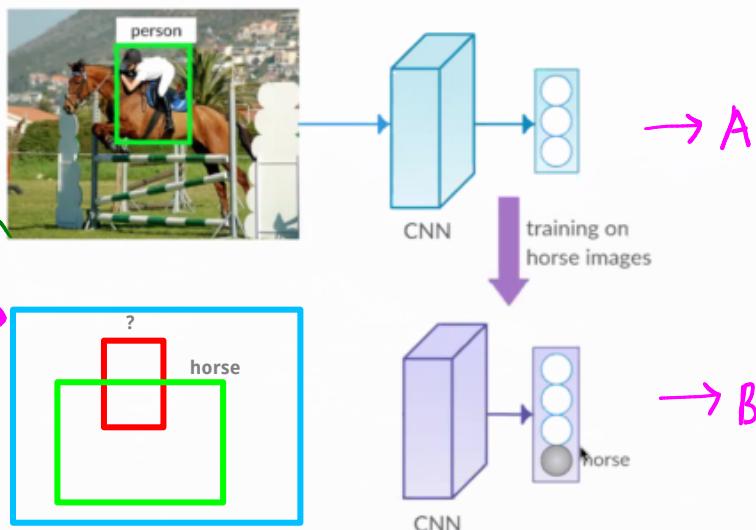
- Finetune the network using new data only
- Leads to catastrophic forgetting



# Incremental Learning: Computer Vision Task

{To be updated}

Dif. Classification  
Task



21

## How well does network B perform ?

method	Training with the <b>initial</b> set of classes	old	new	all
A(1-10)	65.8	-	-	
+B(11-20)	12.8	64.5	<b>38.7</b>	
A(1-20)	68.4	71.3	<b>69.8</b>	
	Baseline, i.e., training with all the classes			

No guidance for retaining the old classes

[Catastrophic forgetting: McCloskey and Cohen 1989, Ratcliff 1990]

26

\* Not retaining what is learnt earlier. (∴ low accuracy)

## Incremental Learning: The Rules !

Task 1

Task 2

Task 3

...

Task n

- Practically {
- Learn one task after the other
  - Without storing (**many**) data from previous tasks
  - Without memory footprint growing (**significantly**) over time
  - Without (**completely**) forgetting old tasks

# What else will we see today?

- Flavour of different approaches:
  1. **Regularization based:** LwF, EBLL, EWC, SI, MAS, IMM, ...
  2. Rehearsal / Replay: iCaRL, DGR, GEM, ...
  3. Architecture based: PackNet, progressive nets , HAT, ...
- More than classification?
- Takeaways

KA: Incremental Learning

30

## Regularization-based Models

- When training a new task,
  - add a regularization term to the loss
  - i.e., term to penalize catastrophic forgetting
- R1: data-focused methods
- R2: model/prior-focused methods

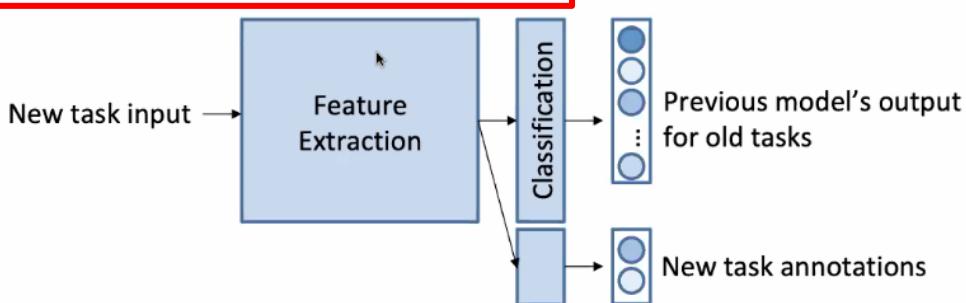


KA: Incremental Learning

Slide credit: T. Tuytelaars

## Data-focused Regularization: Learning without Forgetting

- Knowledge distillation loss
  - i.e., preservation of responses



# Data-focused Regularization: Learning without Forgetting



Simple method; good results for related tasks



Poor results for unrelated tasks

?

Need to store the old model

[Li & Hoiem 2016]

KA: Incremental Learning

33

## Model-focused Regularization

- Penalize changes to ‘important’ parameters

$$\mathcal{L}(\theta) = \underbrace{\mathcal{L}_B(\theta^n)}_{\text{Loss on new task(s)}} + \alpha \sum_k \lambda_k (\theta_k^n - \theta_k^{n-1})^2$$

↑  
Regularization

Different variants possible for  
“importance” and regularization

◀ ▶ ↻

KA: Incremental Learning

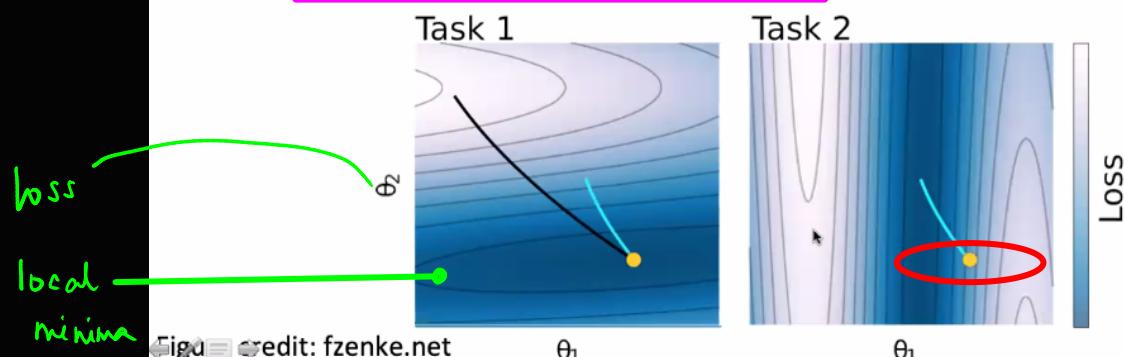
34

## Model-focused Regularization

- Elastic weight consolidation [Kirkpatrick et al., 2017]

- Indiv. penalty for each previous task
- Fisher information matrix for  $\lambda$

→ To read!



Figures credit: fzenke.net

35

# Model-focused Regularization

- Elastic weight consolidation [Kirkpatrick et al., 2017]
    - Indiv. penalty for each previous task  $\sum_k \sum_{i < n} \lambda_k^{n-i} (\theta_k^n - \theta_k^{n-i})^2$
    - Fisher information matrix for  $\lambda$
- 

Figure from paper

KA: Incremental Learning

36

# Model-focused Regularization

- Elastic weight consolidation [Kirkpatrick et al., 2017]
  - Indiv. penalty for each previous task  $\sum_k \sum_{i < n} \lambda_k^{n-i} (\theta_k^n - \theta_k^{n-i})^2$
  - Fisher information matrix for  $\lambda$



Agnostic to architecture; Good results empirically



Only valid locally



Need to store importance weights

KA: Incremental Learning

37

# Model-focused Regularization

- Memory aware synapses [Aljundi et al., 2018]
  - Considers only the previous task  $\sum_k \lambda_k (\theta_k^n - \theta_k^{n-1})^2$
  - Change in gradients for  $\lambda$

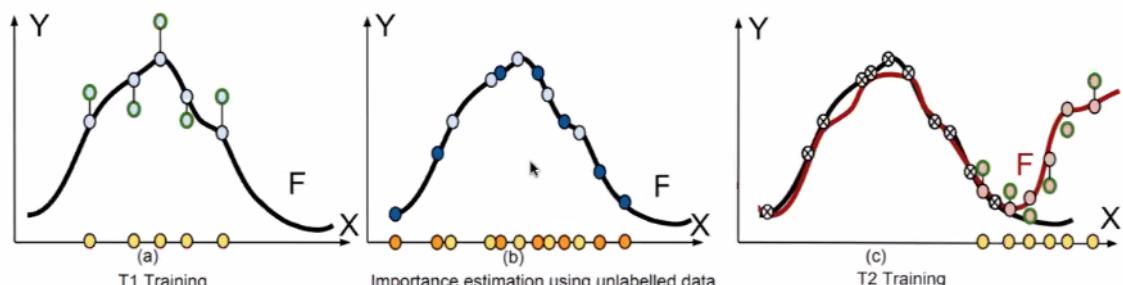


Figure from paper

KA: Incremental Learning

38

# Model-focused Regularization

- Memory aware synapses [Aljundi et al., 2018]

- Considers only the previous task  $\sum_k \lambda_k (\theta_k^n - \theta_k^{n-1})^2$
- Change in gradients for  $\lambda$



Agnostic to architecture; Leverages data & output



Only valid locally



Need to store importance weights

Other examples of  
model focused  
Regularization  
Other than EWS &  
MAS

① Synaptic Intelligence  
② Moment Matching  
③ PathNet.

## What else will we see today?

- Flavour of different approaches:
  1. Regularization based: LwF, EBLL, EWC, SI, MAS, IMM, ...
  2. **Rehearsal / Replay**: iCaRL, DGR, GEM, ...
  3. Architecture based: PackNet, progressive nets , HAT, ...
- More than classification?
- Takeaways

## Rehearsal / Replay-based methods

- Store a couple of examples from previous tasks
- Or produce samples from a generative model
- But
  - How many?
  - How to select them?
  - How to use them?

# iCaRL: Incremental classifier and representation learning

- Selects samples that are closest to the feature mean of each class
- Knowledge distillation loss [Hinton et al.'14]
- Clever use of available memory (see the following)

[Rebuffi et al. 2017]

KA: Incremental Learning

43

→ To  
Study:  
To be  
updated  
iCaRL  
Algorithm \*

## iCaRL: Incremental classifier and representation learning

---

**Algorithm** iCaRL INCREMENTALTRAIN

---

**input**  $X^s, \dots, X^t$  // training examples in per-class sets  
**input**  $K$  // memory size  
**require**  $\Theta$  // current model parameters  
**require**  $\mathcal{P} = (P_1, \dots, P_{s-1})$  // current exemplar sets

$\Theta \leftarrow \text{UPDATEREPRESENTATION}(X^s, \dots, X^t; \mathcal{P}, \Theta)$   
 $m \leftarrow K/t$  // number of exemplars per class  
**for**  $y = 1, \dots, s-1$  **do**  
     $P_y \leftarrow \text{REDUCEEXEMPLARSET}(P_y, m)$   
**end for**  
**for**  $y = s, \dots, t$  **do**  
     $P_y \leftarrow \text{CONSTRUCTEXEMPLARSET}(X_y, m, \Theta)$   
**end for**  
 $\mathcal{P} \leftarrow (P_1, \dots, P_t)$  // new exemplar sets

---

Split the problem into:

- learning features, and then
- using NCM classifier

[Rebuffi et al. 2017]

KA: Incremental Learning

44

## iCaRL: Incremental classifier and representation learning

---

**Algorithm** iCaRL CLASSIFY

---

**input**  $x$  // image to be classified  
**require**  $\mathcal{P} = (P_1, \dots, P_t)$  // class exemplar sets  
**require**  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^d$  // feature map

**for**  $y = 1, \dots, t$  **do**  
     $\mu_y \leftarrow \frac{1}{|P_y|} \sum_{p \in P_y} \varphi(p)$  // mean-of-exemplars  
**end for**  
     $y^* \leftarrow \underset{y=1, \dots, t}{\operatorname{argmin}} \|\varphi(x) - \mu_y\|$  // nearest prototype  
**output** class label  $y^*$

---

[Rebuffi et al. 2017]

KA: Incremental Learning

45

# iCaRL: Incremental classifier and representation learning [Rebuffi et al.'17]

**Algorithm** iCaRL UPDATEREPRESENTATION

**input**  $X^s, \dots, X^t$  // training images of classes  $s, \dots, t$

**require**  $\mathcal{P} = (P_1, \dots, P_{s-1})$  // exemplar sets

**require**  $\Theta$  // current model parameters

// form combined training set:

$$\mathcal{D} \leftarrow \bigcup_{y=s, \dots, t} \{(x, y) : x \in X^y\} \cup \bigcup_{y=1, \dots, s-1} \{(x, y) : x \in P^y\}$$

// store network outputs with pre-update parameters:

**for**  $y = 1, \dots, s-1$  **do**

$$q_i^y \leftarrow g_y(x_i) \quad \text{for all } (x_i, \cdot) \in \mathcal{D}$$

**end for**

run network training (e.g. BackProp) with loss function

$$\ell(\Theta) = - \sum_{(x_i, y_i) \in \mathcal{D}} \left[ \sum_{y=s}^t \delta_{y=y_i} \log g_y(x_i) + \delta_{y \neq y_i} \log(1 - g_y(x_i)) \right. \\ \left. + \sum_{y=1}^{s-1} q_i^y \log g_y(x_i) + (1 - q_i^y) \log(1 - g_y(x_i)) \right]$$

Classification loss

← Distillation loss:  
Comparing old vs new

that consists of *classification* and *distillation* terms.

[Rebuffi et al. 2017]

KA: Incremental Learning

46

## iCaRL: Incremental classifier and representation learning



Clever use of available memory



Potential issues with storing data, e.g., privacy



Limited by the memory capacity (the more the better)

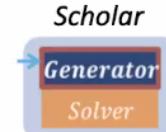
[Rebuffi et al. 2017]

KA: Incremental Learning

47

## Deep Generative Replay

- The model “Scholar” is composed of:
  - a generator + a solver (classifier)
- The generator and the solver are updated in every incremental step



\* Just like  
GAN's  
generator &  
discrimi-  
nator

[Shin et al. 2017]

Figure from the paper

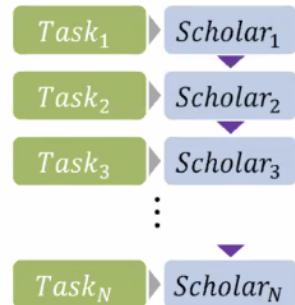
KA: Incremental Learning

49

# Deep Generative Replay

Training procedure:

- At task  $t$ , we train a new Scholar
  - with data from the task  $t$ , and
  - data generated by the previously trained Scholar at task  $t-1$



[Shin et al. 2017]

Figure from the paper

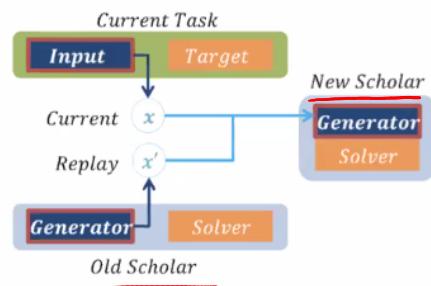
KA: Incremental Learning

Slide courtesy: A. Massenet 50

## Deep Generative Replay

Training procedure (Generator):

- With data from task  $t$ , and
- data generated by the previously trained Scholar for task  $t-1$



[Shin et al. 2017]

Figure from the paper

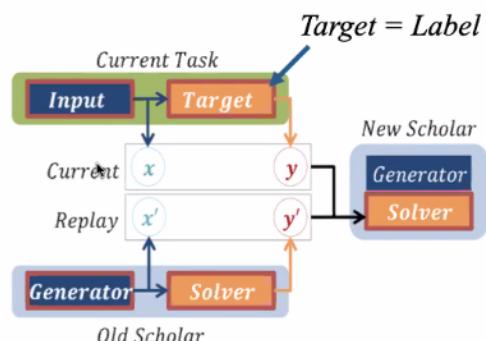
KA: Incremental Learning

Slide courtesy: A. Massenet 51

## Deep Generative Replay

Training procedure (Solver):

- With data from task  $t$ , and
- Data from generator and solver of the previously trained Scholar for task  $t-1$



[Shin et al. 2017]

Figure from the paper

KA: Incremental Learning

Slide courtesy: A. Massenet 52

# Deep Generative Replay



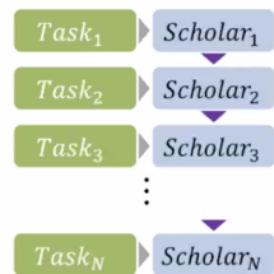
Avoids memory issues



Accumulation of errors



No control over the class of the generated samples



[Shin et al. 2017]

Figure from the paper

KA: Incremental Learning

Slide courtesy: A. Massenet 53

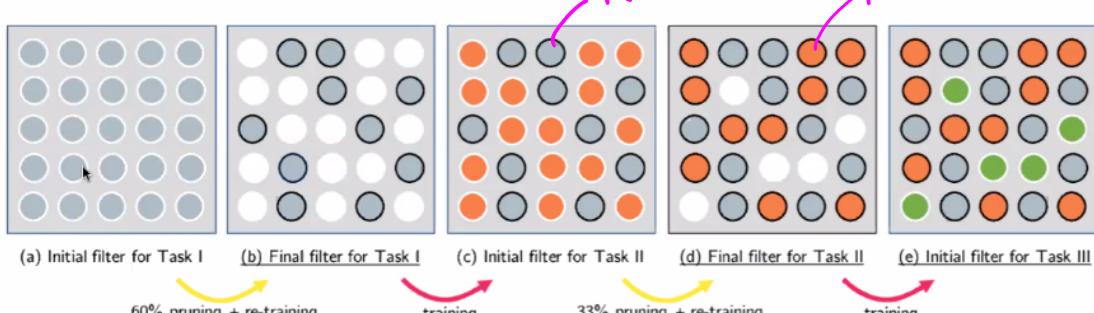
## What else will we see today?

- Flavour of different approaches:
  1. Regularization based: LwF, EBLL, EWC, SI, MAS, IMM, ...
  2. Rehearsal / Replay: iCaRL, DGR, GEM, ...
  3. **Architecture based:** PackNet, progressive nets , HAT, ...
- More than classification?
- Takeaways

KA: Incremental Learning

54

## Architecture-based



But needs to know prior the total no. of tasks.

PackNet [Mallya & Lazebnik'17]

Figure from the paper

KA: Incremental Learning

55

# Architecture-based



Fixed memory consumption



Needs the total number of tasks



Avoids forgetting

PackNet [Mallya & Lazebnik'17]

KA: Incremental Learning

56

## A Comparative Analysis

- TinyImagenet: small, balanced, class-incremental
- iNaturalist: large-scale, unbalanced, task-incremental

	Tiny Imagenet	iNaturalist
Tasks	10	10
Classes per task	20	5 to 314
Training data per task	8k	0.6k to 66k
Validation data per task	1k	0.1k to 9k
Task Constitution	random class selection	supercategory



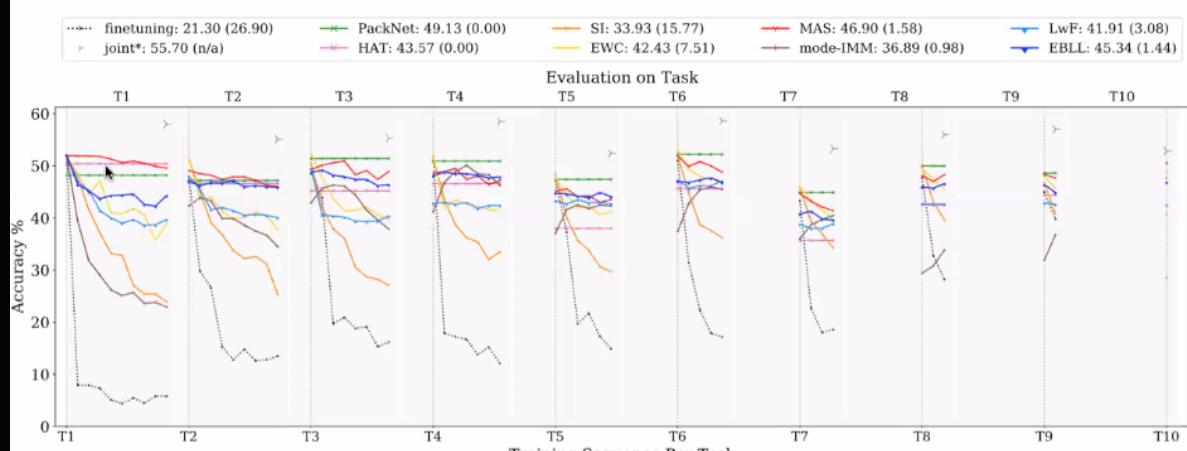
- Fair way of setting hyperparameters (stability-plasticity tradeoff)

Lange et al., 2020]

KA: Incremental Learning

57

## Comparative Evaluation (TinyImagenet)



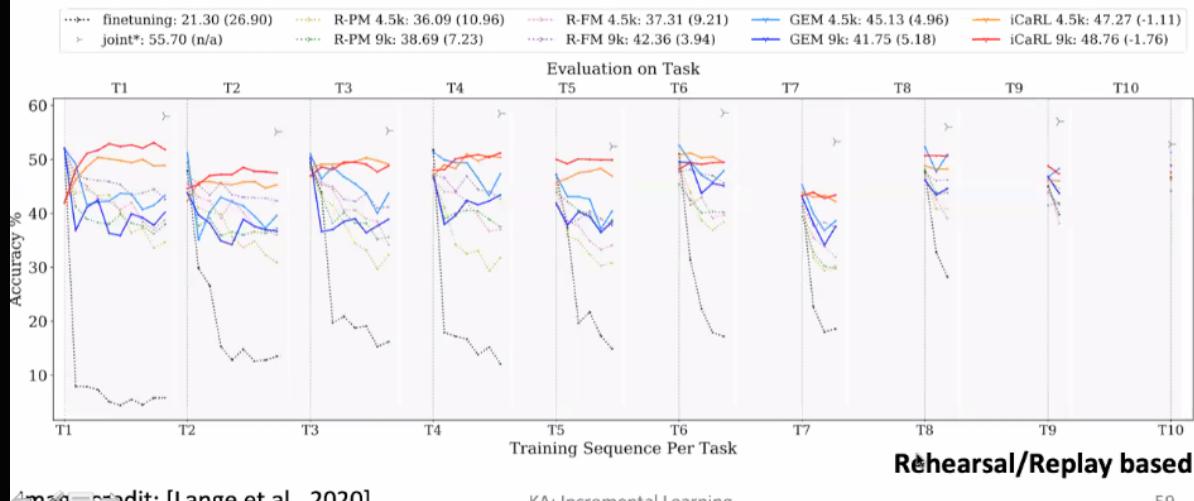
Regularization & Architecture based

← → credit: [Lange et al., 2020]

KA: Incremental Learning

58

# Comparative Evaluation (TinyImageNet)



← mag edit: [Lange et al., 2020]

KA: Incremental Learning

59

## \* General Trends

- • Rehearsal/replay based methods only pay off when storing significant amount of exemplars
- • PackNet results in no-forgetting and produces top results
- • MAS more robust than EWC

← mag edit: T. Tuytelaars

KA: Incremental Learning

60

## Summary

- Flavour of different approaches:
  1. Regularization based: LwF, EBLL, EWC, SI, MAS, IMM, ...
  2. Rehearsal / Replay: iCaRL, DGR, GEM, ...
  3. Architecture based: PackNet, progressive nets , HAT, ...
- More than classification?
- Takeaways

← mag

KA: Incremental Learning

82

## Looking to the future

- Desiderata
  - Constant memory
  - Task agnostic: Some recent advances [Rao et al., NeurIPS'19]
  - Forgetting gracefully
  - Datasets

“I don’t like datasets, it’s more a problem than a solution” – heard at ICCV 2019