

# CVPR 2020 Tutorial

## Towards Annotation-Efficient Learning

### Incremental Learning

Karteek Alahari

Inria, France

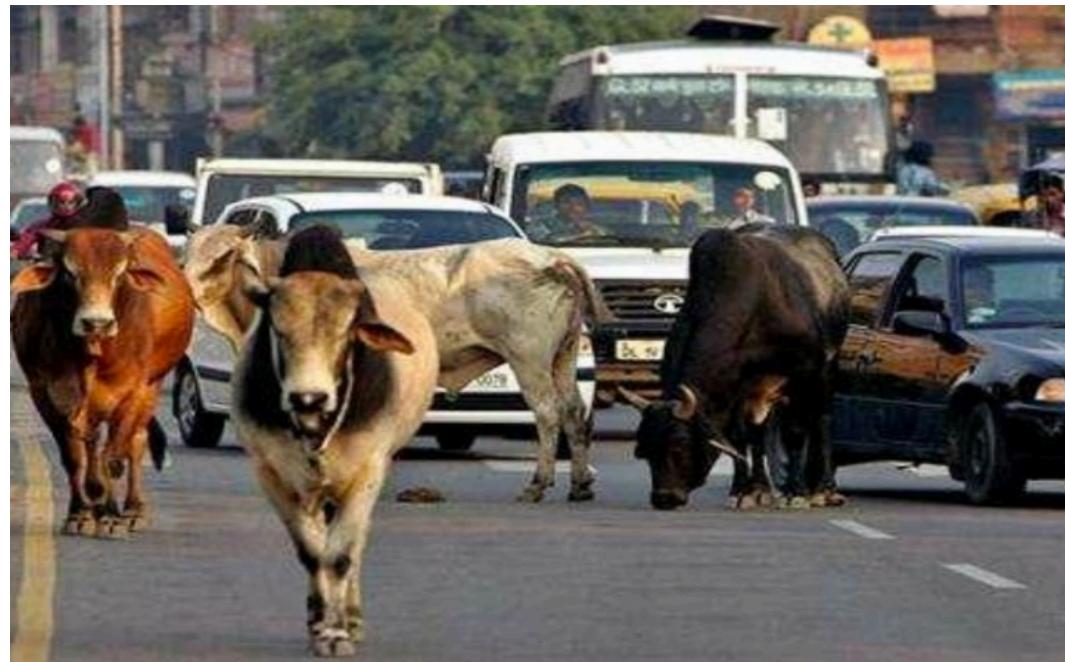


# Incremental Learning ?

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

# An Incremental Learning Scenario

- Growing up in India



KA: Incremental Learning

# An Incremental Learning Scenario

- And then during travels



Robert Berdan ©

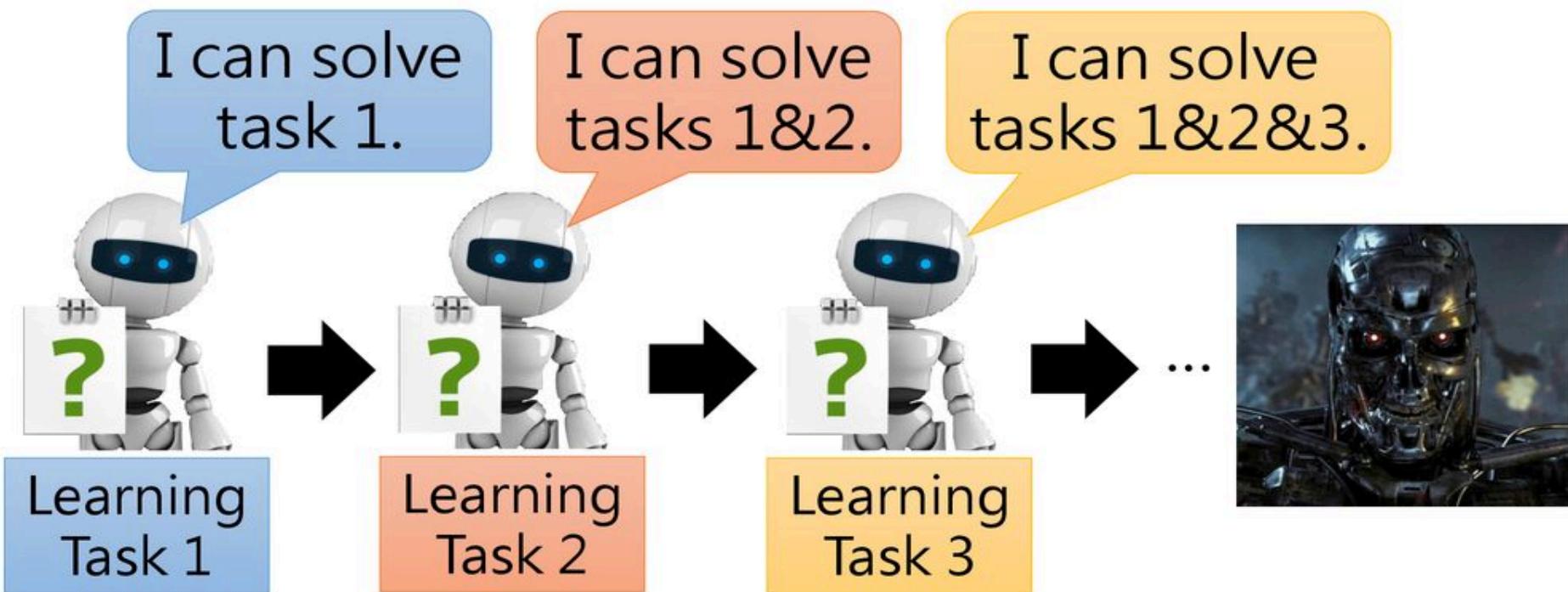
KA: Incremental Learning

# An Incremental Learning Scenario

- And then during travels



KA: Incremental Learning



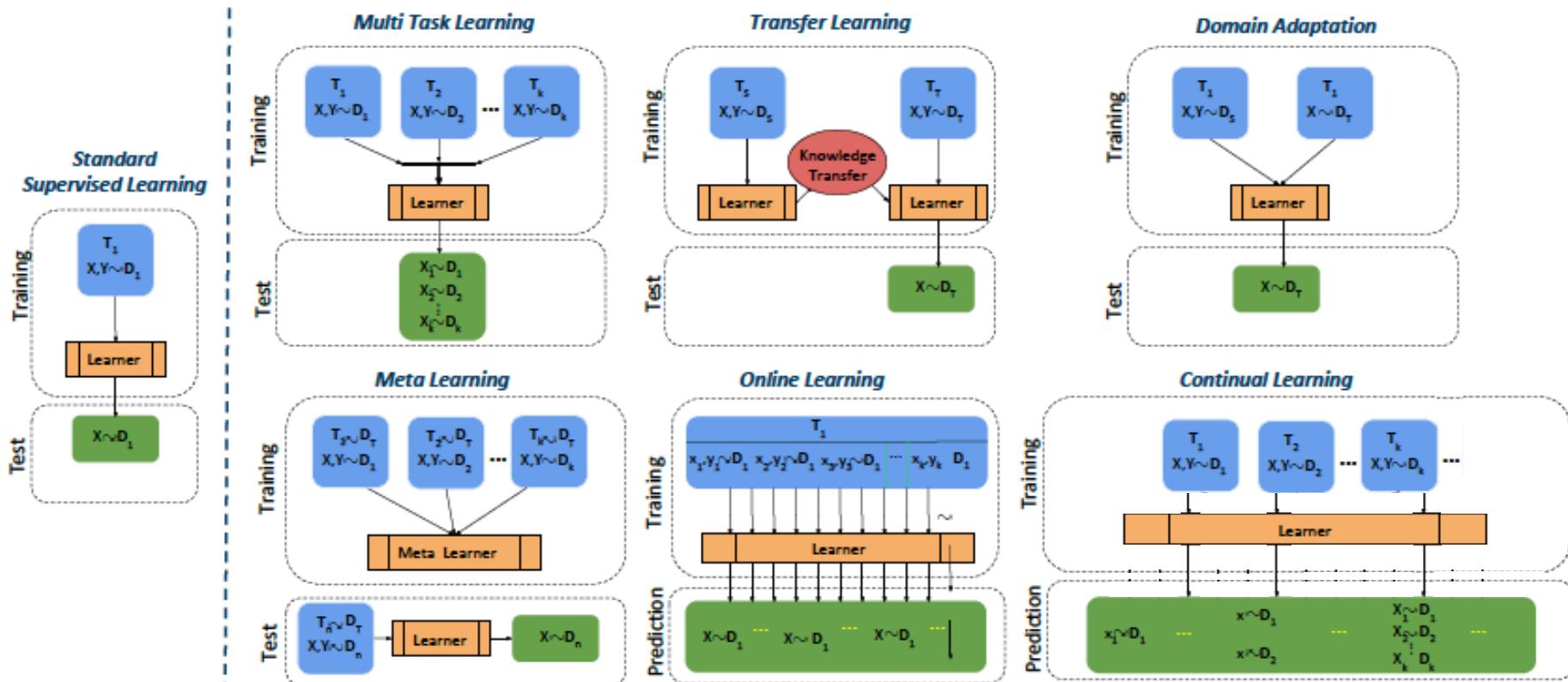
KA: Incremental Learning

Slide credit: Hung-yi Lee<sup>6</sup>

# Standard Machine Learning

TRAIN – VALIDATION – TEST

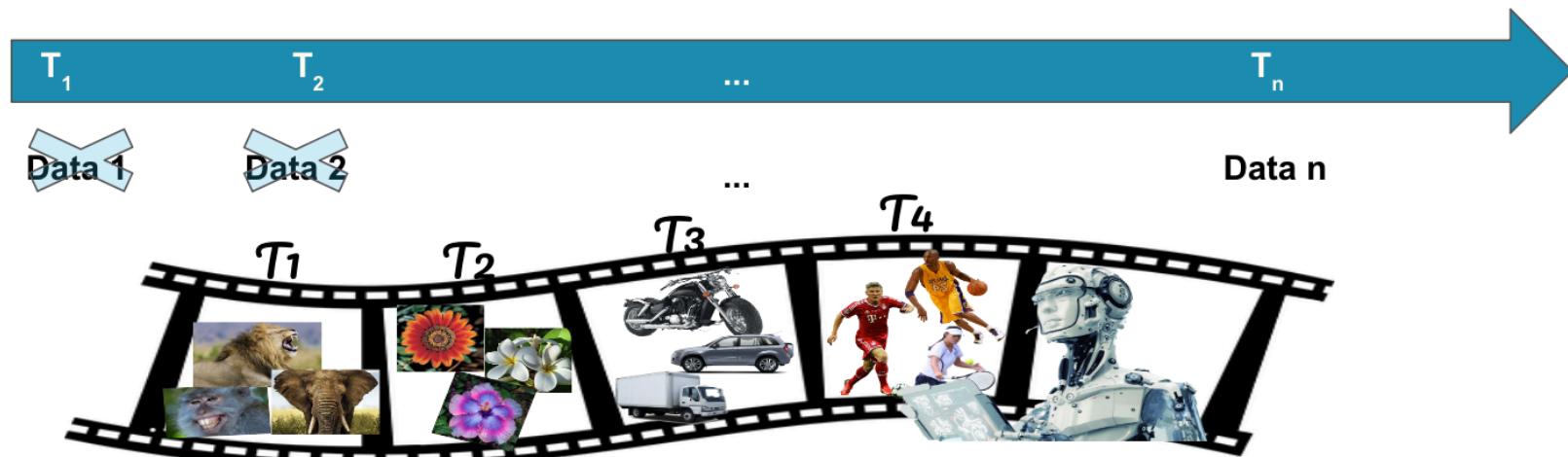
- All sampled from the same distribution
- > benchmarks and academic datasets 😊
- > real-world systems ⚡
- > embodied learning ⚡



KA: Incremental Learning

Slide credit: T. Tuytelaars

# Incremental Learning Setup



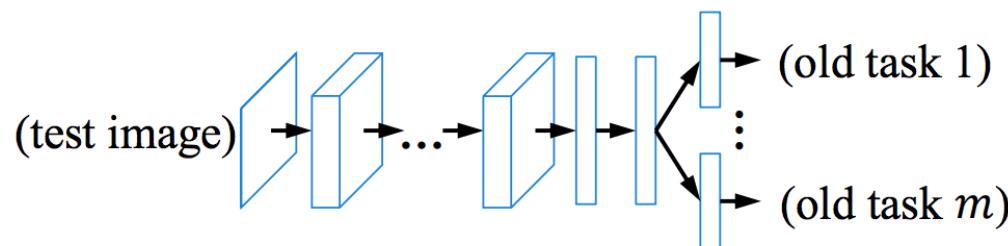
- Task-incremental learning
- Class-incremental learning
- Domain-incremental learning

# 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

# Brute Force Solution (non-incremental)

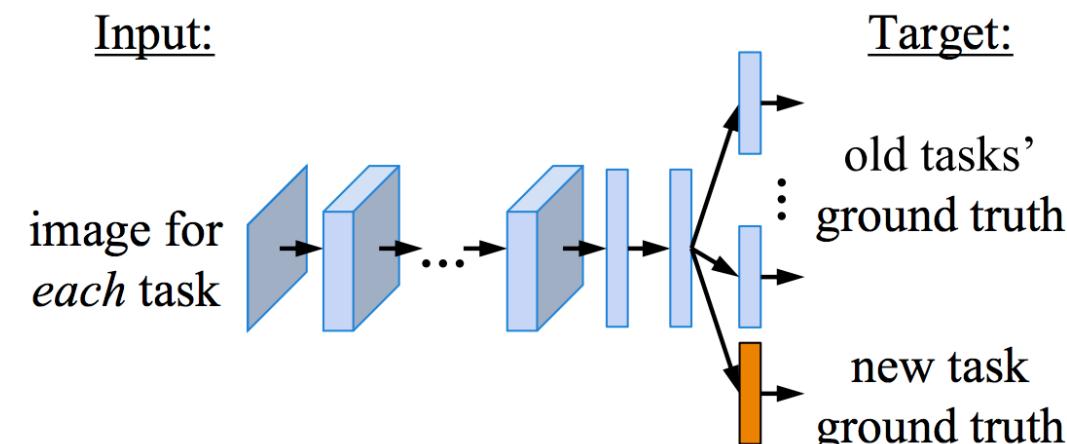
Original model



Joint training

“golden” baseline

- random initialize + train
- fine-tune
- unchanged



KA: Incremental Learning

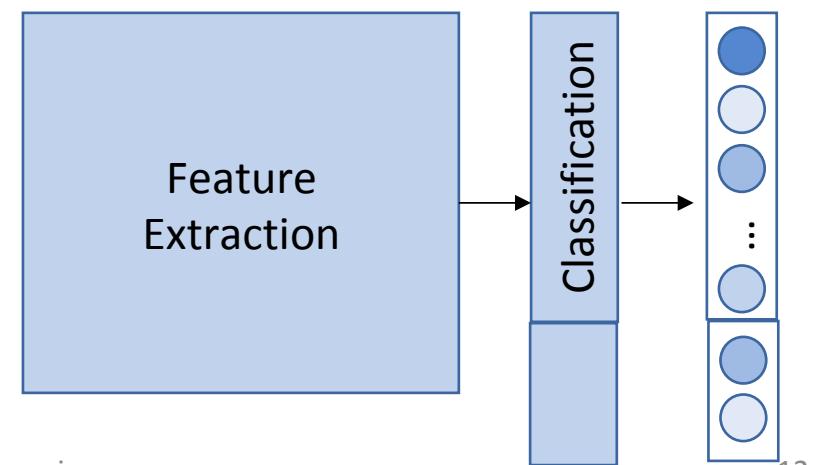
Figures from [Li and Hoiem 2016]<sup>11</sup>

# Brute Force Solution (non-incremental)

Retrain full model with both old and new data

- Computationally expensive

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

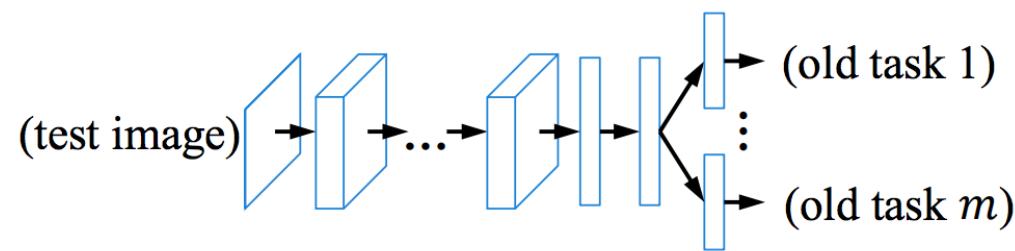


# Why not brute force ?

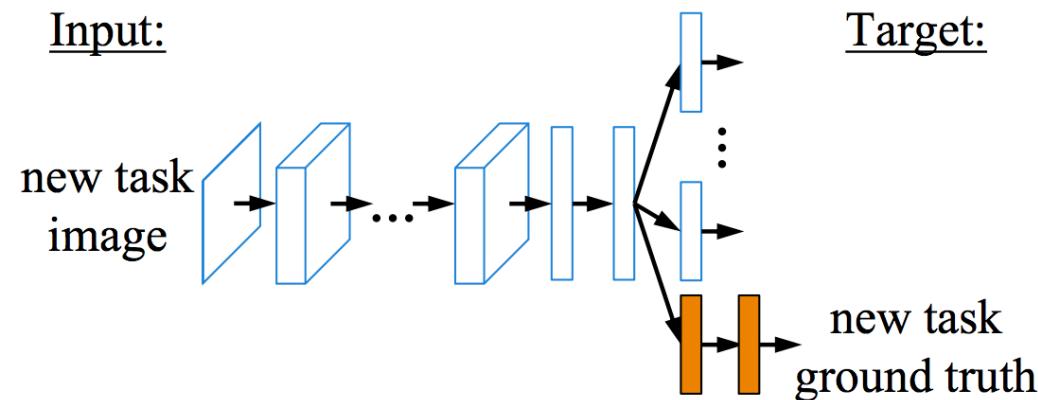
- No access to all the data
- Can not store all the data
- Access to only a previously learned model, e.g., trained by others

# Naïve Solution 1

Original model



Feature extraction



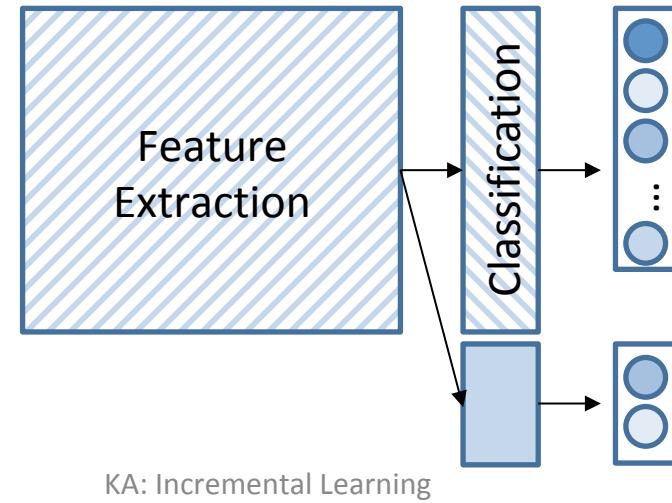
- random initialize + train
- fine-tune
- unchanged

KA: Incremental Learning

Figures from [Li and Hoiem 2016]

# Naïve Solution 1

- Finetune only last layer using new data only
  - Leads to **suboptimal results**

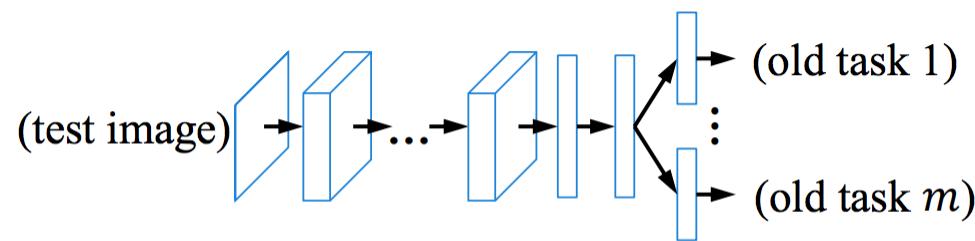


KA: Incremental Learning

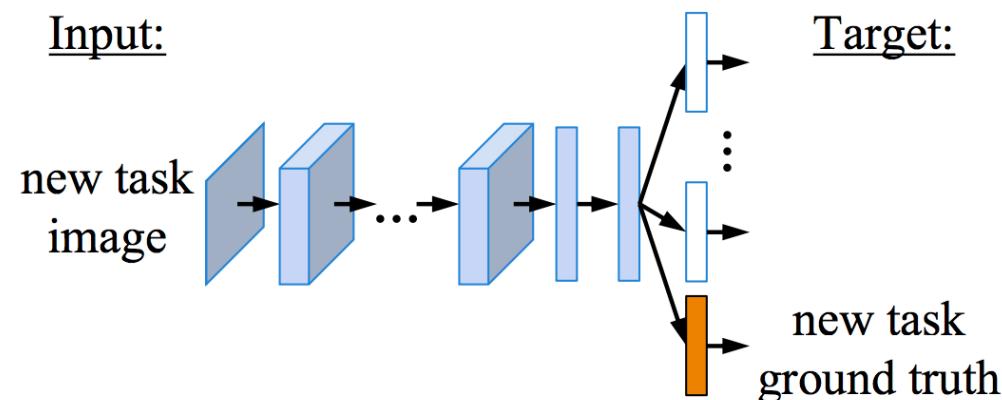
Slide credit: T. Tuytelaars<sup>15</sup>

# Naïve Solution 2

Original model



Fine tuning



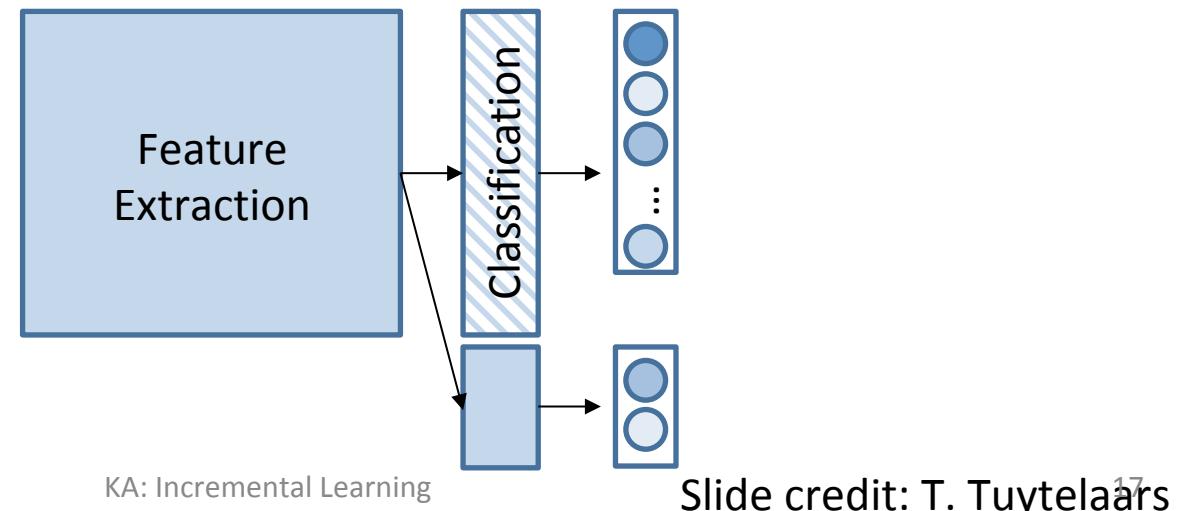
- random initialize + train
- fine-tune
- unchanged

KA: Incremental Learning

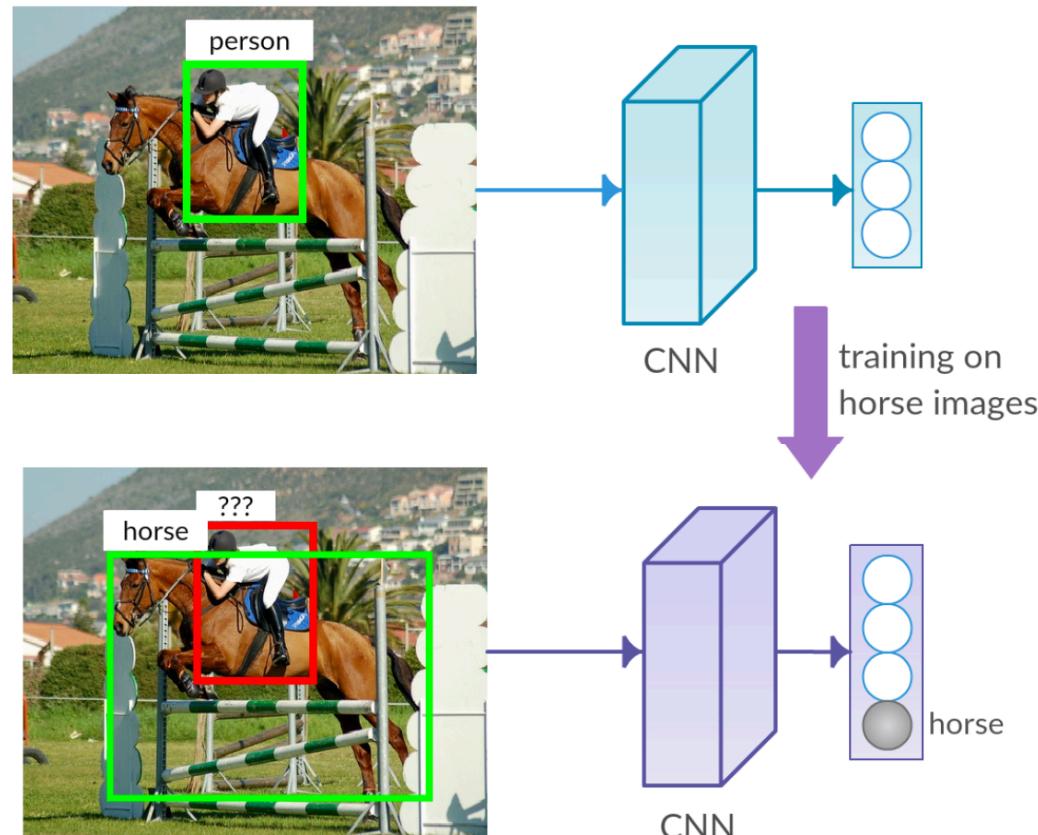
Figures from [Li and Hoiem 2016]

# Naïve Solution 2

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



# Incremental Learning: Computer Vision Task



A

B

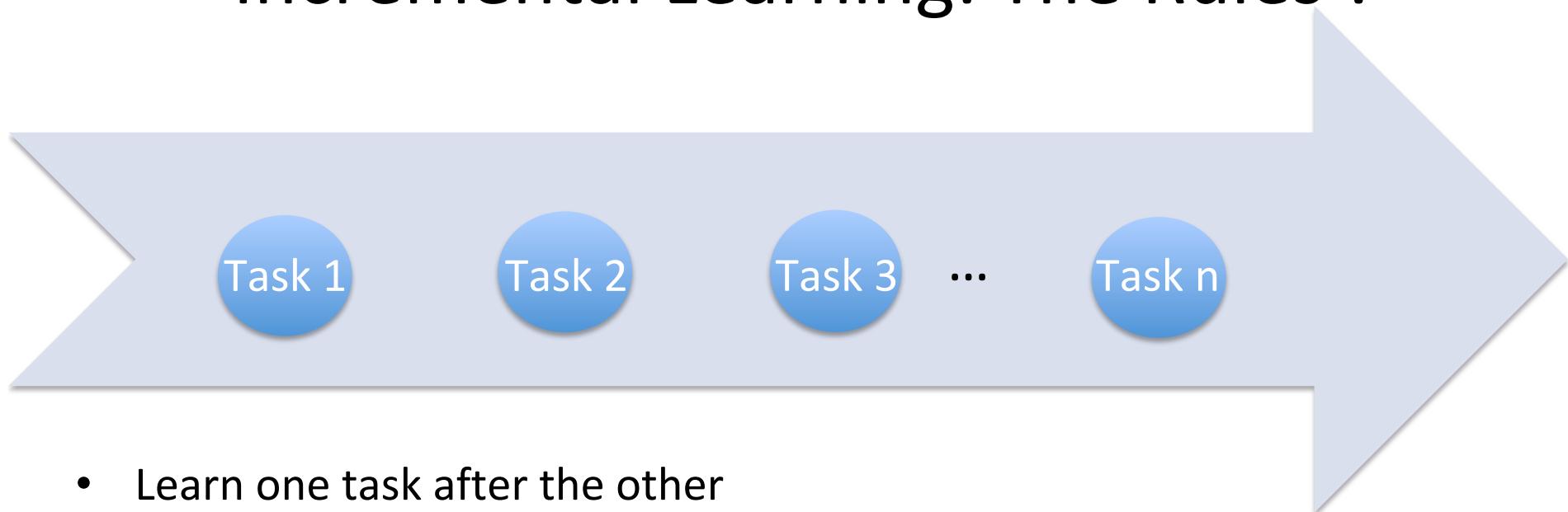
# How well does network B perform ?

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

No guidance for retaining the old classes

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

# Incremental Learning: The Rules !



- 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

# What else will we see today?

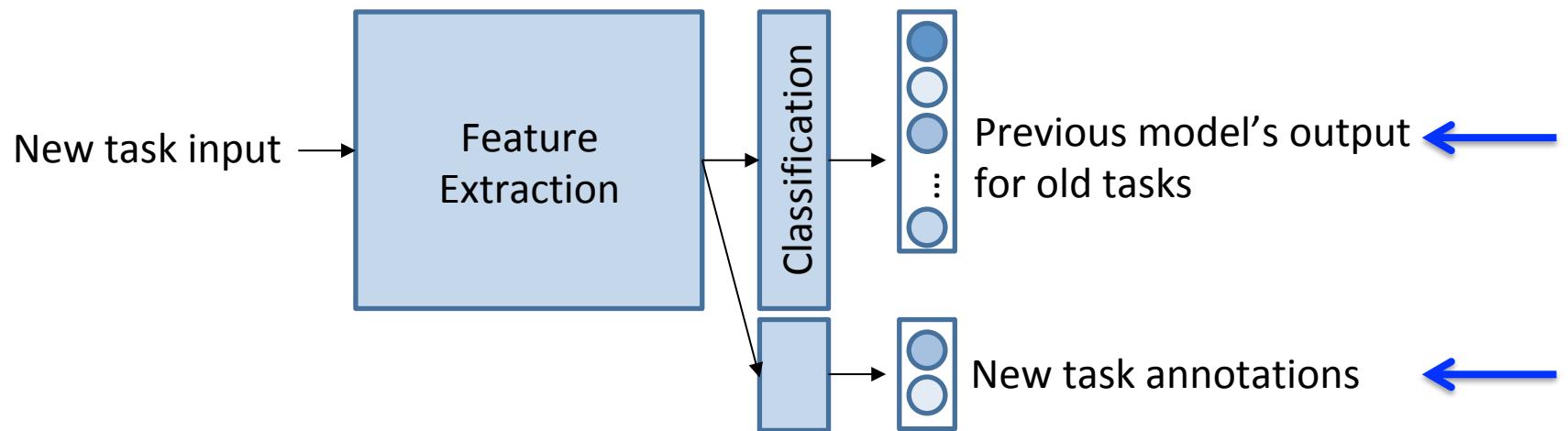
- Flavour of different approaches:
  1. **Regularization based**: LwF, EBLL, EWC, SI, MAS, IMM, ...
  2. Rehearsal / Replay: iCaRL, DGR, GEM, ...
  3. Architecture based: PackNet, progressive nets , HAT, ...
- More than classification?
- Takeaways

# 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

# Data-focused Regularization: Learning without Forgetting

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



[Li & Hoiem 2016]

KA: Incremental Learning

Slide credit: T. Tuytelaars<sup>34</sup>

# Data-focused Regularization: Learning without Forgetting



Simple method; good results for related tasks



Poor results for unrelated tasks



Need to store the old model

# 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**

# Model-focused Regularization

- Elastic weight consolidation [Kirkpatrick et al., 2017]
  - Indiv. penalty for each previous task
  - Fisher information matrix for  $\lambda$

$$\sum_k \sum_{i < n} \lambda_k^{n-i} (\theta_k^n - \theta_k^{n-i})^2$$

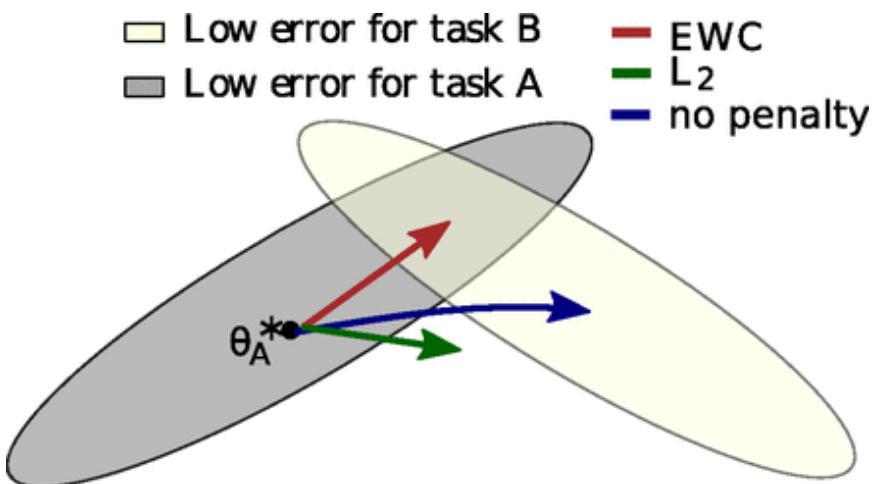


Figure from paper

KA: Incremental Learning

27

# 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                   |

# 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  $\lambda$

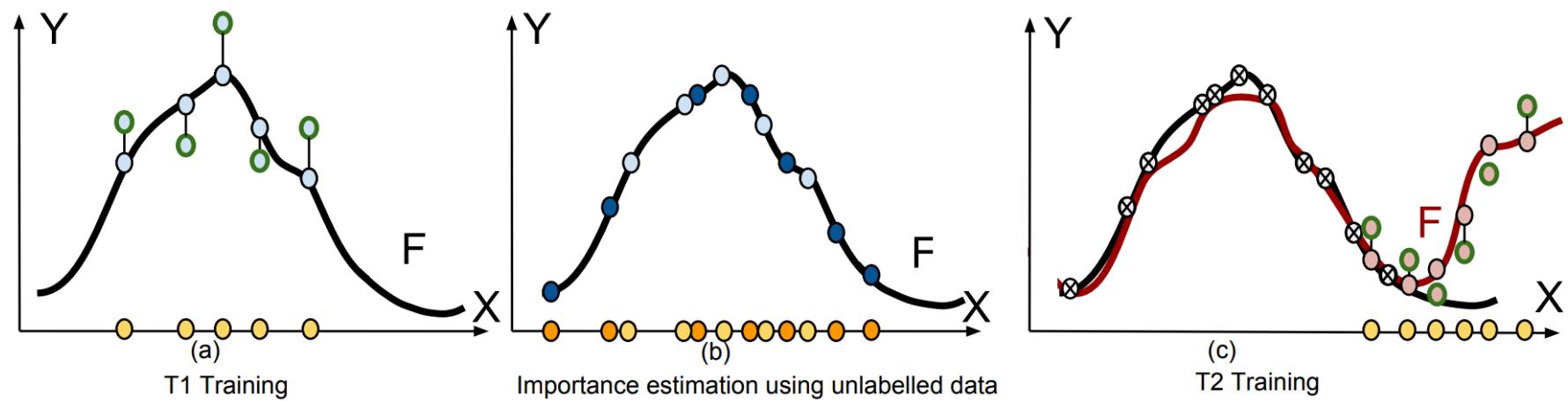


Figure from paper

KA: Incremental Learning

# 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

# Model-focused Regularization

- Two examples
  - Elastic weight consolidation [Kirkpatrick et al., 2017]
  - Memory aware synapses [Aljundi et al., 2018]
- Other alternatives
  - Path Integral / Synaptic Intelligence: large changes during training [Zenke et al., 2017]
  - Moment matching [Lee et al., 2017]
  - Pathnet [Fernando et al., 2017]
  - ...

# 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)

# iCaRL: Incremental classifier and representation learning

Split the problem into:

- learning features, and then
- using NCM classifier

---

**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
```

---



[Rebuffi et al. 2017]

# iCaRL: Incremental classifier and representation learning

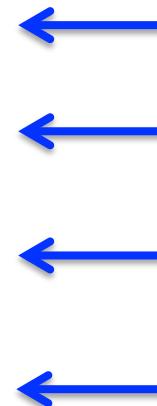
---

**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  
     $y^* \leftarrow \operatorname{argmin}_{y=1, \dots, t} \|\varphi(x) - \mu_y\|$  // nearest prototype  
output class label  $y^*$ 
```

---



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

---

**Algorithm** iCaRL UPDATEREPRESENTATION

---

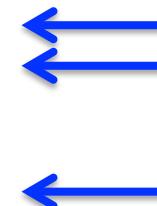
**input**  $X^s, \dots, X^t$  // training images of classes  $s, \dots, t$   
**require**  $\mathcal{P} = (P_1, \dots, P_{s-1})$  // exemplar sets  
**require**  $\Theta$  // current model parameters  
// form combined training set:

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

// store network outputs with pre-update parameters:  
**for**  $y = 1, \dots, s-1$  **do**  
 $q_i^y \leftarrow g_y(x_i)$  for all  $(x_i, \cdot) \in \mathcal{D}$   
**end for**  
run network training (e.g. BackProp) with loss function

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

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



Classification loss



Distillation loss:  
Comparing old vs new

# iCaRL: Incremental classifier and representation learning



Clever use of available memory



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



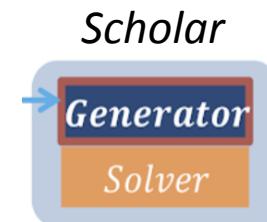
Limited by the memory capacity (the more the better)

# 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

# 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



[Shin et al. 2017]

Figure from the paper

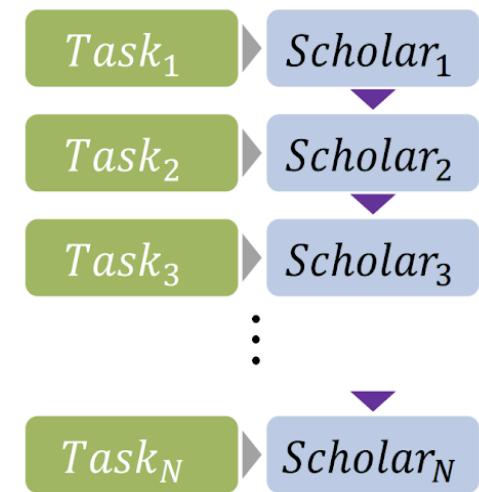
KA: Incremental Learning

40

# 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

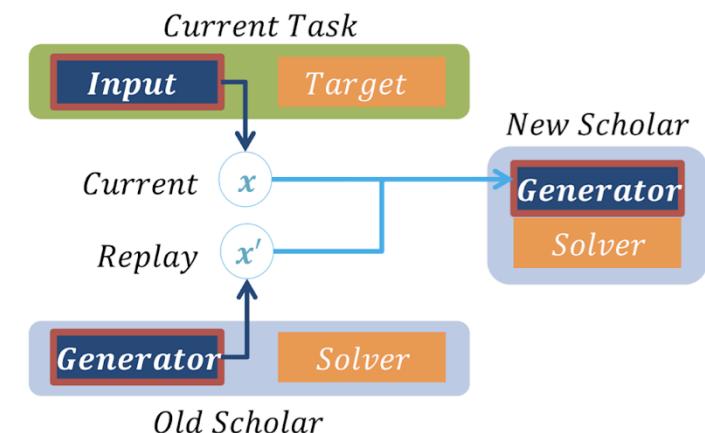
KA: Incremental Learning

Slide courtesy: A. Massenet 41

# 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 42

# Deep Generative Replay

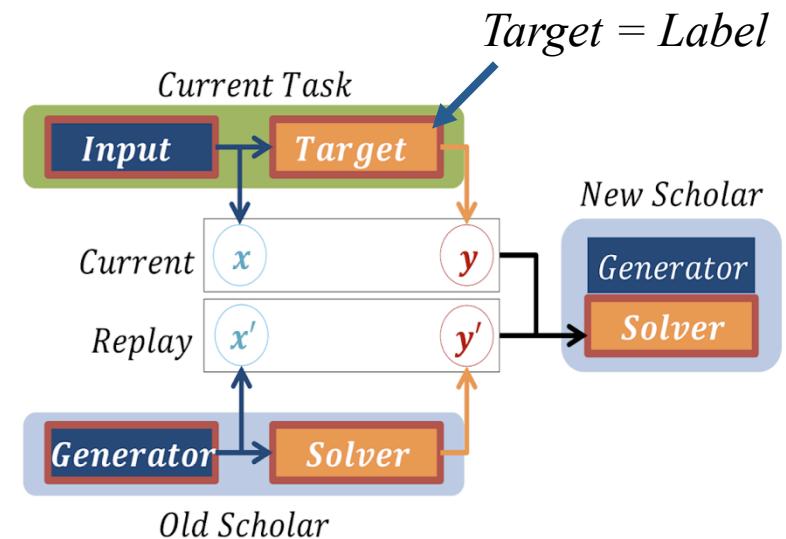
Training procedure (Solver):

- 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 43

# Deep Generative Replay



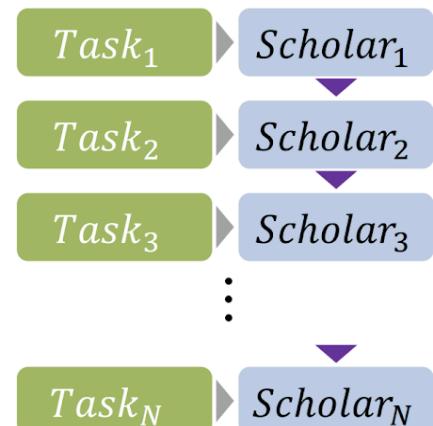
Avoids memory issues



Accumulation of errors



No control over the class of the generated samples



[Shin et al. 2017]

Figure from the paper

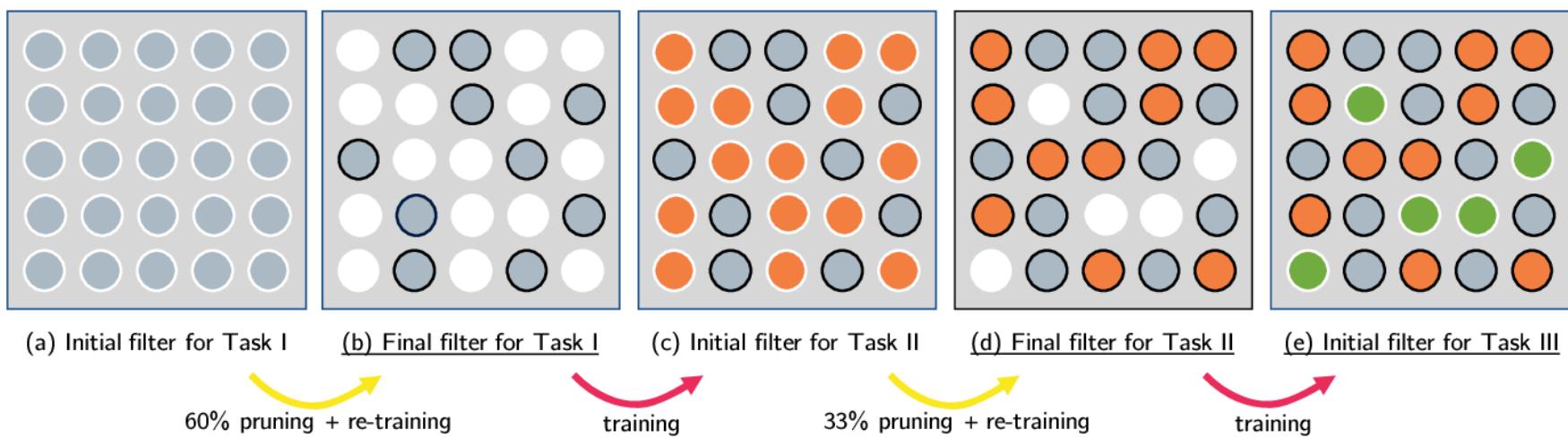
KA: Incremental Learning

Slide courtesy: A. Massenet 44

# 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

# Architecture-based



PackNet [Mallya & Lazebnik'17]

Figure from the paper

KA: Incremental Learning

46

# Architecture-based



Fixed memory consumption



Needs the total number of tasks



Avoids forgetting

PackNet [Mallya & Lazebnik'17]

KA: Incremental Learning

47

# 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)



# Comparative Evaluation (TinylImagenet)

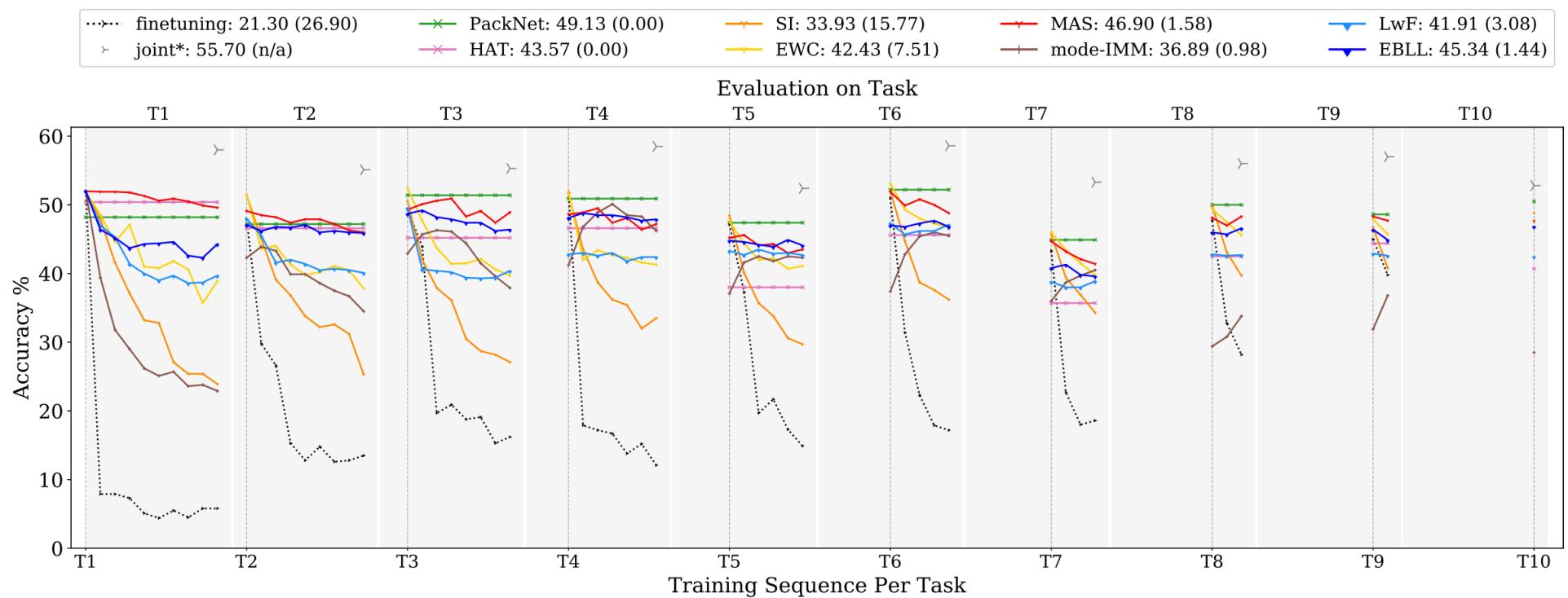


Image credit: [Lange et al., 2020]

KA: Incremental Learning

49

# Comparative Evaluation (TinyImagenet)

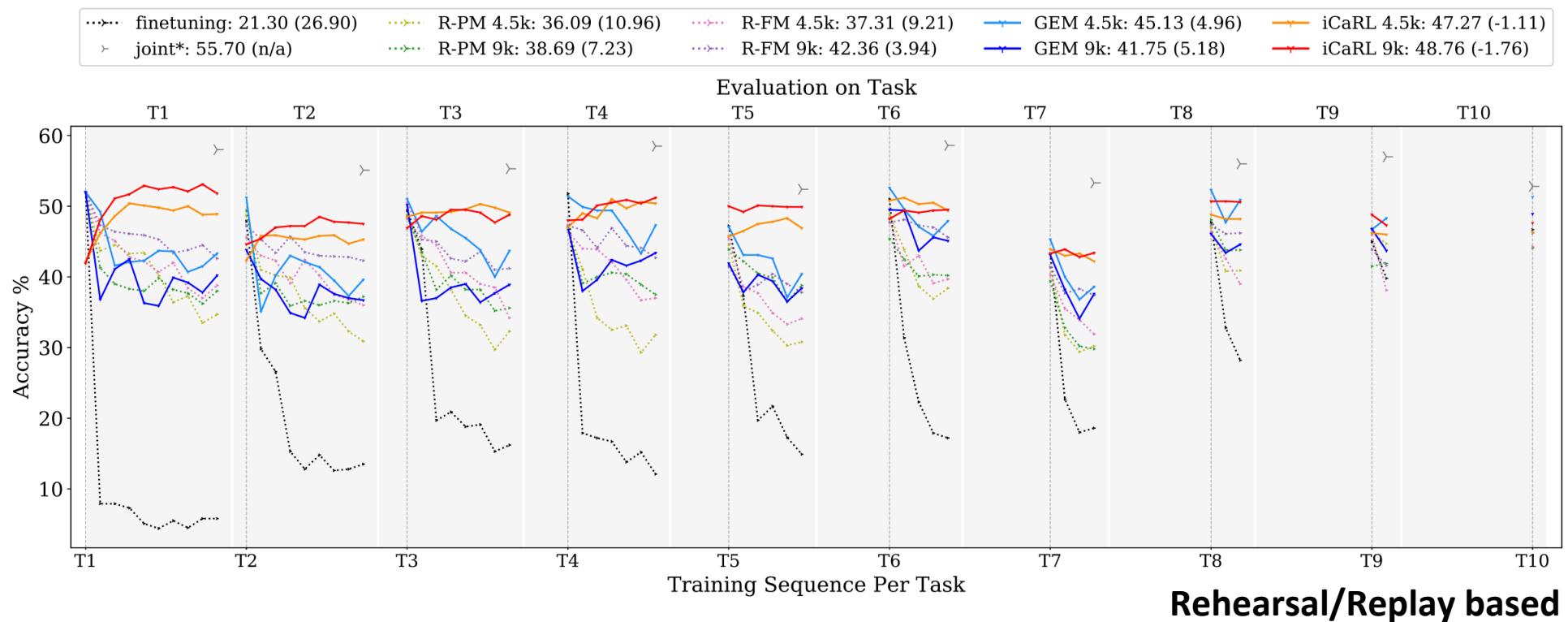


Image credit: [Lange et al., 2020]

KA: Incremental Learning

50

# 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

# What kind of model should I use ?

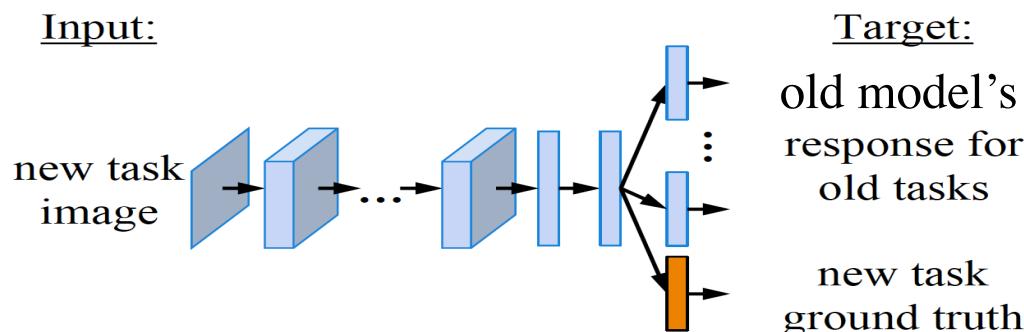
- Larger models give more capacity (but: overfitting)
  - Wide is better than deep
- 
- Regularization may interfere with incremental learning
  - Dropout usually better than weight decay

# 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

# Mitigate Catastrophic Forgetting

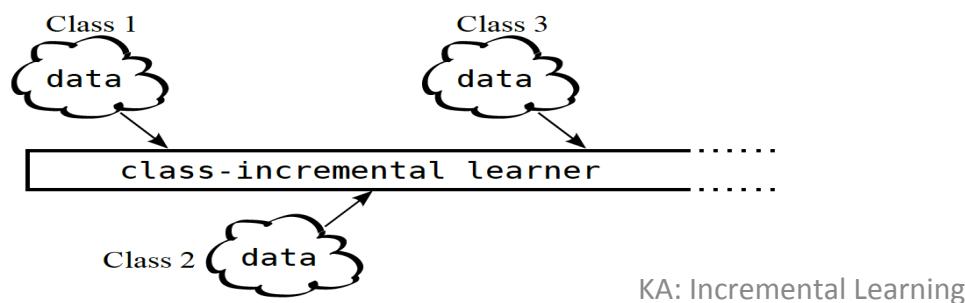
- Learning without forgetting [Li and Hoiem 2016]



Tasks defined on a new dataset

Focus on image classification  
(rare co-occurrence of old and new)

- iCaRL [Rebuffi et al. 2017]



Decouple classifier and feature learning

Rely on a subset of the old data

# Mitigate Catastrophic Forgetting

- Elastic weight consolidation [Kirkpatrick et al., 2017]
  - Selectively slowing down learning on weights
    - Limited to specific settings
    - Focus on image classification  
(rare co-occurrence of old and new)
- Other attempts, e.g., [Aljundi et al., 2018, Jung et al., 2016, Mallya and Lazebnik, 2017, Risin et al., 2014, Rusu et al., 2016]

# Mitigate Catastrophic Forgetting

- Elastic weight consolidation [Kirkpatrick et al., 2017]
  - Selectively slowing down learning on weights
    - Limited to specific settings
- Other attempts, e.g., [Jung et al., 2016, Mallya and Lazebnik, 2017, Risin et al., 2014, Rusu et al., 2016]

**Lack of methods for  
incremental learning of object detectors**

# An approach

- Incremental Learning of Object Detectors without Catastrophic Forgetting [Shmelkov et al., 2017]

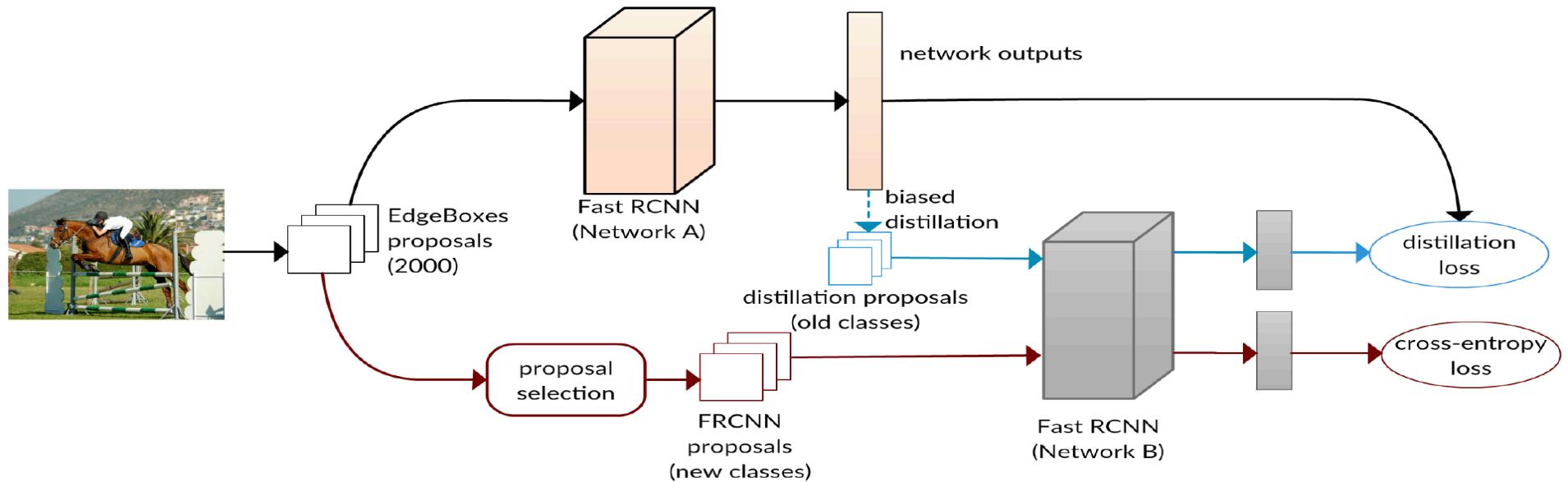


# Our Approach

- Builds on knowledge distillation [Hinton et al. '14]
- Knowledge distillation originally for:
  - learning a simpler net equivalent to a complex one, and
  - produces a differently structured network
- We distill the “old” net when learning new classes

# Our Approach

- Copy of network (**A**) to evaluate proposals & memorize outputs



- New component (**B**) to learn the new classes
- Learn with a combination of losses

KA: Incremental Learning

[Fast RCNN: Girshick 2015]

# Our Approach

- Minimize

$$\mathcal{L} = \mathcal{L}_{\text{rcnn}} + \lambda \mathcal{L}_{\text{dist}}$$

$$\mathcal{L}_{\text{rcnn}}(\mathbf{p}, k^*, t, t^*) = - \underbrace{\log \mathbf{p}_{k^*}}_{\text{Loss over classes}} + \underbrace{[k^* \geq 1] R(t - t^*)}_{\text{Localization loss}}$$

$$\mathcal{L}_{\text{dist}}(y_A, t_A, y_B, t_B) = \frac{1}{N|C_A|} \sum \left[ \underbrace{(\bar{y}_A - \bar{y}_B)^2 + (t_A - t_B)^2}_{\text{Comparing the two net components}} \right]$$

# Our Approach

- Minimize

$$\mathcal{L} = \mathcal{L}_{\text{rcnn}} + \lambda \mathcal{L}_{\text{dist}}$$

$$\mathcal{L}_{\text{rcnn}}(\mathbf{p}, k^*, t, t^*) = - \underbrace{\log \mathbf{p}_{k^*}}_{\text{Loss over classes}} + \underbrace{[k^* \geq 1] R(t - t^*)}_{\text{Localization loss}}$$

$$\mathcal{L}_{\text{dist}}(y_A, t_A, y_B, t_B) = \frac{1}{N|C_A|} \sum \left[ \underbrace{(\bar{y}_A - \bar{y}_B)^2}_{\text{Logit responses}} + \underbrace{(t_A - t_B)^2}_{\text{Localization}} \right]$$

# Experimental setup

- Datasets: PASCAL VOC 2007 and MS COCO
- A(n-m): initial training on classes from n to m (old)
- B(n-m): incremental training on classes n to m (new) added all at once
- B(n)(m): incremental training for class n, then m, added one after another
- Unbiased distillation: distillation proposals are selected completely randomly

## Addition of 10 classes

method	old	new	all
A(1-10)	65.8	-	-
+B(11-20), no distillation	12.8	64.5	38.7
+B(11-20) with distillation	63.2	63.1	63.1
A(1-20) (baseline)	68.4	71.3	69.8

## Addition of 10 classes

method	old	new	all
A(1-10)	65.8	-	-
+B(11-20), no distillation	12.8	64.5	38.7
+B(11-20) with distillation	<b>63.2</b>	63.1	63.1
+B(11-20) with EWC	<b>31.6</b>	61.0	46.3
A(1-20) (baseline)	68.4	71.3	69.8

# Addition of 5 classes

- VOC 2007 validation set  
method

	old	new	all
A(1-15)	70.5	-	-
+B(16-20) with distillation	<b>68.4</b>	58.4	65.9
+B(16)(17)...(20) with distillation	<b>66.0</b>	51.6	62.4
A(1-20) (baseline)	70.9	66.7	69.8

# Addition of 5 classes

- VOC 2007 validation set  
method

	old	new	all
A(1-15)	70.5	-	-
+B(16-20) with distillation	68.4	58.4	65.9
+B(16)(17)...(20) with distillation	<b>66.0</b>	51.6	62.4
+B(16)(17)...(20) w unbiased distillation	<b>45.8</b>	46.5	46.0
A(1-20) (baseline)	70.9	66.7	69.8

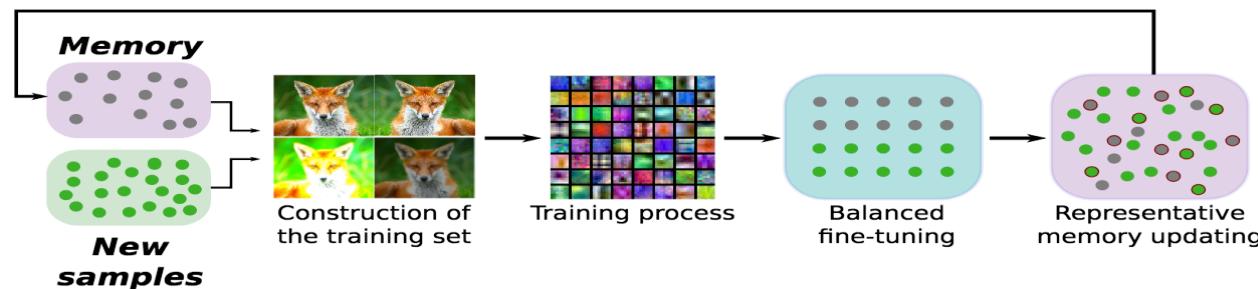
## Addition of 40 classes

- COCO minival (first 5000 validation images)

method	mAP@.5	mAP@[.5, .95]
A(1-40)+B(41-80)	37.4	21.3
A(1-80) (baseline)	38.1	22.6

# Intermediate Summary

- No catastrophic forgetting with
  - distillation loss
  - in particular, biased distillation
- At ECCV 2018: Balanced fine-tuning



# 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**

# 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

# 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

See also: Workshop on 14<sup>th</sup> June  
Continual Learning in Computer Vision