

Overcoming Catastrophic Forgetting with Unlabeled Data in the Wild

Presenters: Nikhil Kannan, Ying Fan

1 Introduction:

1.1 Catastrophic Forgetting:

- Goal of class-incremental learning is to learn a model that performs well on previous and new tasks without task boundaries. But it suffers from catastrophic forgetting.
- Training Neural Networks on new tasks causes it to forget information learned from previously trained tasks, degrading model performance on earlier tasks.
- Primary reason for catastrophic forgetting is limited resources for scalability.

1.2 Class Incremental Learning Setting

- $(x, y) \in \mathbb{D}$, T is a supervised task mapping $x \rightarrow y$
- For task T_t , corresponding dataset is \mathbb{D}_t and coreset is $\mathbb{D}^{\text{cor}}_{t-1} \subseteq \mathbb{D}_{t-1} \cup \mathbb{D}^{\text{cor}}_{t-2}$ contains representative data of previous tasks $T_{1:(t-1)} = \{T_1, \dots, T_t\}$. For task T_t corresponding labeled training data used is represented as $\mathbb{D}_t^{\text{trn}} = \mathbb{D}_t \cup \mathbb{D}^{\text{cor}}_{t-1}$.
- $M_t = \{\theta, \phi_{1:t}\}$ is a set of learnable parameters of a model where θ indicates shared task parameters and $\phi_{1:t} = \{\phi_1, \dots, \phi_t\}$ are task specific parameters.

2. Local distillation & Global distillation

2.1 Local distillation:

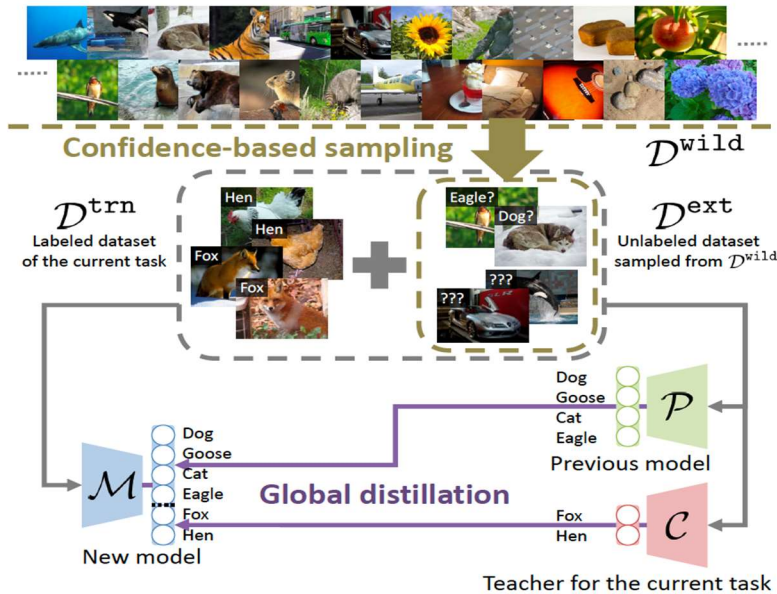
- Train the model M_t by minimizing the classification loss: $L_{\text{cls}}(\theta, \phi_{1:t}; \mathbb{D}_t^{\text{trn}})$.
- In the class incremental learning setting, the limited capacity of coreset causes the model to suffer from catastrophic forgetting. To overcome this issue, utilize previously trained model M_{t-1} , that contains knowledge of previous tasks to generate soft labels: Optimize $\sum_{s=1}^{t-1} L_{\text{dst}}(\theta, \phi_s; P_t, \mathbb{D}_t)$, where $P_t = \{\theta^P, \phi_{1:(t-1)}^P\} = M_{t-1}$ is a previous trained model

- Then minimize the joint objective: $L_{\text{cls}}(\theta, \phi_{1:t}; \mathbb{D}_t^{\text{trn}}) + \sum_{s=1}^{t-1} L_{\text{dst}}(\theta, \phi_s; P_t, \mathbb{D}_t)$
- Solving the above optimization problem is called local knowledge distillation. Transfers the knowledge within each of the tasks. The issue with local knowledge distillation is that is defined in a task-wise manner and misses the knowledge about discriminating between classes in different tasks.

2.2 Global distillation:

- Distill the knowledge of reference models globally by minimizing the following loss: $L_{\text{dst}}(\theta, \phi_{1:(t-1)}; P_t, \mathbb{D}_t^{\text{trn}} \cup \mathbb{D}_t^{\text{ext}})$
- Learning using the above function causes bias, since P_t does not have knowledge regarding the current task, hence performance on the current task is degraded. So introduce teacher model $C_t = \{\theta^C, \phi^C_t\}$ specialized to learn the current task $T_t: L_{\text{dst}}(\theta, \phi_t; C_t, \mathbb{D}_t^{\text{trn}} \cup \mathbb{D}_t^{\text{ext}})$, where teacher model C_t is trained by minimizing $L_{\text{cls}}(\theta^C, \phi^C_t; \mathbb{D}_t)$
- P_t learns to perform tasks $T_{1:(t-1)}$ and C_t learns to perform the current task T_t , but knowledge distillation between $T_{1:(t-1)}$ and T_t is not captured by the either of the reference models. Define Q_t , an ensemble of reference models P_t and C_t : ensemble $Q_t: L_{\text{dst}}(\theta, \phi_{1:t}; Q_t, \mathbb{D}_t^{\text{ext}})$
- The global distillation model learns by optimizing the following loss:

$$L_{\text{cls}}(\theta, \phi_{1:t}; \mathbb{D}_t^{\text{trn}}) + L_{\text{dst}}(\theta, \phi_{1:(t-1)}; P_t, \mathbb{D}_t^{\text{trn}} \cup \mathbb{D}_t^{\text{ext}}) + L_{\text{dst}}(\theta, \phi_t; C_t, \mathbb{D}_t^{\text{trn}} \cup \mathbb{D}_t^{\text{ext}}) + L_{\text{dst}}(\theta, \phi_{1:t}; Q_t, \mathbb{D}_t^{\text{ext}})$$



3. Fine-Tuning and Normalization

3.1 Normalization:

- Since the amount of data from the previous tasks is smaller than the current task, model prediction is biased towards the current task. To remove the bias, fine tune the model after the training phase by scaling the computed gradient from the data with label k .

- $w_D^{(k)} = \frac{1}{|\{(x, y) \in \mathbb{D} \mid y=k\}|}$, scaling the gradient is similar to feeding data multiple times (data weighting). Normalizing weights by multiplying them with $|\mathbb{D}|/|T|$ to balance the dataset \mathbb{D}

3.2 Fine-tuning:

- Fine-tune task-specific ($\phi_{1:t}$) using data weighting to remove any bias from training data and to equally weigh training data for all tasks. Also, fine-tuning shared parameters (θ) is not required since it already contains relevant information from all training data.

- Loss Weight: balance the contribution of each loss by the relative size of each task learned in the loss; $w^L = \frac{|T|}{|T_{1:t}|}$

4. 3-step Learning Algorithm

Learning strategy has three steps:

- Training C_t specialized for learning the current task T_t
- Training M_t through global knowledge distillation of reference models P_t, Q_t, C_t
- Fine-tuning model parameters using data weighting.

Algorithm 1 3-step learning with GD.

```

1:  $t = 1$ 
2: while true do
3:   Input: previous model  $\mathcal{P}_t = \mathcal{M}_{t-1}$ , coresets  $\mathcal{D}_{t-1}^{\text{cor}}$ ,
      training dataset  $\mathcal{D}_t$ , unlabeled data stream  $\mathcal{D}_t^{\text{wild}}$ 
4:   Output: new coresets  $\mathcal{D}_t^{\text{cor}}$ , model  $\mathcal{M}_t = \{\theta, \phi_{1:t}\}$ 
5:    $\mathcal{D}_t^{\text{trn}} = \mathcal{D}_t \cup \mathcal{D}_{t-1}^{\text{cor}}$ 
6:    $N_C = |\mathcal{D}_{t-1}^{\text{cor}}|, N_D = |\mathcal{D}_t^{\text{trn}}|$ 
7:   Sample  $\mathcal{D}_t^{\text{ext}}$  from  $\mathcal{D}_t^{\text{wild}}$  using Algorithm 2
8:   Train  $C_t$  by minimizing Eq. (12)
9:   if  $t > 1$  then
10:    Train  $\mathcal{M}_t$  by minimizing Eq. (9)
11:    Fine-tune  $\phi_{1:t}$  by minimizing Eq. (9),
      with data weighting in Eq. (10)
12:   else
13:     $\mathcal{M}_t = C_t$ 
14:   end if
15:   Randomly sample  $\mathcal{D}_t^{\text{cor}} \subseteq \mathcal{D}_t^{\text{trn}}$  such that
       $|\{(x, y) \in \mathcal{D}_t^{\text{cor}} | y = k\}| = N_C / |\mathcal{T}_{1:t}|$  for  $k \in \mathcal{T}_{1:t}$ 
16:    $t = t + 1$ 
17: end while

```

5. Sampling External Dataset

5.1 The main issues with using unlabeled data in knowledge distillation.

- Training is computationally expensive
- Most of the unlabeled data might be irrelevant to the tasks the model learns

The paper proposes a sampling method to sample an external dataset from large stream of unlabeled data that benefits knowledge distillation:

5.2 Confidence Calibration

Sampling external data that is expected to be in previous tasks is desirable, since it alleviates catastrophic forgetting. Neural Nets tend to be highly overconfident as they produce prediction with high confidence for OOD data. To achieve confidence calibrated outputs, model learns from certain amount of OOD data and data from previous tasks:

- For the model to produce confidence calibrated outputs, following confidence loss function is considered: $L_{\text{cnf}}(\theta, \phi; \mathbb{D}) = \frac{1}{|\mathbb{D}||T|} \sum_{x \in \mathbb{D}} \sum_{y \in T} [-\log p(y|x; \theta, \phi)]$
- During 3-step learning, training C_t has no reference model hence it learns from confidence loss. By optimizing on confidence loss, model learns to produce predictions with low confidence for OOD data.
- C_t learns by optimizing $L_{\text{cls}}(\theta^C, \phi^C_t; \mathbb{D}_t) + L_{\text{cnf}}(\theta^C, \phi^C_t; \mathbb{D}_{t-1}^{\text{cor}} \cup \mathbb{D}_t^{\text{ext}})$

Algorithm 2 Sampling external dataset.

```

1: Input: previous model  $\mathcal{P}_t = \{\theta^{\mathcal{P}}, \phi_{1:(t-1)}^{\mathcal{P}}\}$ ,
   unlabeled data stream  $\mathcal{D}_t^{\text{w}^{\text{IID}}}$ , sample size  $N_D$ ,
   number of unlabeled data to be retrieved  $N_{\text{max}}$ 
2: Output: sampled external dataset  $\mathcal{D}_t^{\text{ext}}$ 
3:  $\mathcal{D}^{\text{prev}} = \emptyset, \mathcal{D}^{\text{OOD}} = \emptyset$ 
4:  $N_{\text{prev}} = 0.3N_D, N_{\text{OOD}} = 0.7N_D$ 
5:  $N(k) \triangleq |\{(x, y, p) \in \mathcal{D}^{\text{prev}} | y = k\}|$ 
6: while  $|\mathcal{D}^{\text{OOD}}| < N_{\text{OOD}}$  do
7:   Get  $x \in \mathcal{D}_t^{\text{w}^{\text{IID}}}$  and update  $\mathcal{D}^{\text{OOD}} = \mathcal{D}^{\text{OOD}} \cup \{x\}$ 
8: end while
9:  $N_{\text{ret}} = N_{\text{OOD}}$ 
10: while  $N_{\text{ret}} < N_{\text{max}}$  do
11:   Get  $x \in \mathcal{D}_t^{\text{w}^{\text{IID}}}$  and compute the prediction of  $\mathcal{P}$ :
      $\hat{p} = \max_y p(y|x; \theta^{\mathcal{P}}, \phi_{1:(t-1)}^{\mathcal{P}})$ ,
      $\hat{y} = \arg \max_y p(y|x; \theta^{\mathcal{P}}, \phi_{1:(t-1)}^{\mathcal{P}})$ 
12:   if  $N(\hat{y}) < N_{\text{prev}}/|\mathcal{T}_{1:(t-1)}|$  then
13:      $\mathcal{D}^{\text{prev}} = \mathcal{D}^{\text{prev}} \cup \{(x, \hat{y}, \hat{p})\}$ 
14:   else
15:     Replace the least probable data in class  $\hat{y}$ :
      $(x', \hat{y}, p') = \arg \min_{(x, y, p) \in \mathcal{D}^{\text{prev}} | y = \hat{y}} p$ 
16:     if  $p' < \hat{p}$  then
17:        $\mathcal{D}^{\text{prev}} = (\mathcal{D}^{\text{prev}} \setminus \{(x', \hat{y}, p')\}) \cup \{(x, \hat{y}, \hat{p})\}$ 
18:     end if
19:   end if
20:    $N_{\text{ret}} = N_{\text{ret}} + 1$ 
21: end while
22: Return  $\mathcal{D}_t^{\text{ext}} = \mathcal{D}^{\text{OOD}} \cup \{x | (x, y, p) \in \mathcal{D}^{\text{prev}}\}$ 

```

6. Related Work

6.1 Continual lifelong learning: **class/task/data** incremental learning

6.2 Methods: model-based and data-based

- Model based: parameters for new tasks are directly constrained to be around that for previous tasks
- Data based: data distribution from previous tasks are used to distill knowledge for later tasks; previous works focus on task-wise local distillation, previous state-of-art: LwF, DR, E2E.

6.3 Proposed method: GD

7. Experiments

7.1 Experimental settings:

- Labeled: CIFAR_100, ImageNet ILSVRC 2012
- Unlabeled: TinyImages, ImageNet2011
- Design tasks: total 100 classes, divide into splits of 5,10,20--task size: 20,10,5
- Hyper parameters: WRN-16-2, coreset size=2000, temperature for smoothing softmax probabilities: 2 for P,C, 1 for Q

7.2 Metrics:

The accuracy of the s-th model at r-th task, $s \geq r$:

$$A_{r,s} = \frac{1}{|\mathcal{D}_r^{\text{test}}|} \sum_{(x,y) \in \mathcal{D}_r^{\text{test}}} \mathbb{I}(\hat{y}(x; \mathcal{M}_s) = y)$$

ACC: weighted combination of accuracy from all tasks and all models:

$$\text{ACC} = \frac{1}{t-1} \sum_{s=2}^t \sum_{r=1}^s \frac{|\mathcal{T}_r|}{|\mathcal{T}_{1:s}|} A_{r,s}.$$

FGT: weighted combination of performance decay:

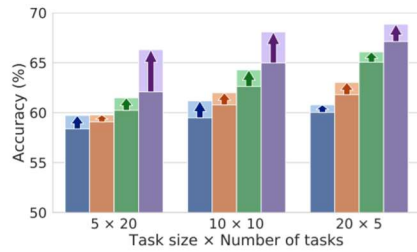
$$\text{FGT} = \frac{1}{t-1} \sum_{s=2}^t \sum_{r=1}^{s-1} \frac{|\mathcal{T}_r|}{|\mathcal{T}_{1:s}|} (A_{r,r} - A_{r,s}),$$

7.3 Results:

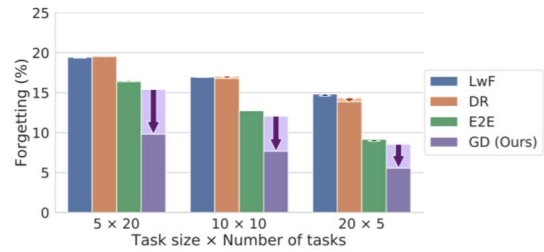
- Overall performance:

Table 1. Comparison of methods on CIFAR-100 and ImageNet. We report the mean and standard deviation of ten trials for CIFAR-100 and nine trials for ImageNet with different random seeds in %. \uparrow (\downarrow) indicates that the higher (lower) number is the better.

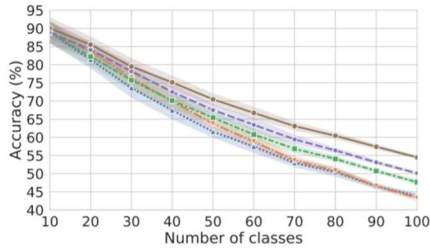
Dataset	CIFAR-100						ImageNet					
	5		10		20		5		10		20	
Metric	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)
Oracle	78.6 \pm 0.9	3.3 \pm 0.2	77.6 \pm 0.8	3.1 \pm 0.2	75.7 \pm 0.7	2.8 \pm 0.2	68.0 \pm 1.7	3.3 \pm 0.2	66.9 \pm 1.6	3.1 \pm 0.3	65.1 \pm 1.2	2.7 \pm 0.2
Baseline	57.4 \pm 1.2	21.0 \pm 0.5	56.8 \pm 1.1	19.7 \pm 0.4	56.0 \pm 1.0	18.0 \pm 0.3	44.2 \pm 1.7	23.6 \pm 0.4	44.1 \pm 1.6	21.5 \pm 0.5	44.7 \pm 1.2	18.4 \pm 0.5
LwF [24]	58.4 \pm 1.3	19.3 \pm 0.5	59.5 \pm 1.2	16.9 \pm 0.4	60.0 \pm 1.0	14.5 \pm 0.4	45.6 \pm 1.9	21.5 \pm 0.4	47.3 \pm 1.5	18.5 \pm 0.5	48.6 \pm 1.2	15.3 \pm 0.6
DR [12]	59.1 \pm 1.4	19.6 \pm 0.5	60.8 \pm 1.2	17.1 \pm 0.4	61.8 \pm 0.9	14.3 \pm 0.4	46.5 \pm 1.6	22.0 \pm 0.5	48.7 \pm 1.6	18.8 \pm 0.5	50.7 \pm 1.2	15.1 \pm 0.5
E2E [3]	60.2 \pm 1.3	16.5 \pm 0.5	62.6 \pm 1.1	12.8 \pm 0.4	65.1 \pm 0.8	8.9 \pm 0.2	47.7 \pm 1.9	17.9 \pm 0.4	50.8 \pm 1.5	13.4 \pm 0.4	53.9 \pm 1.2	8.8 \pm 0.3
GD (Ours)	62.1 \pm 1.2	15.4 \pm 0.4	65.0 \pm 1.1	12.1 \pm 0.3	67.1 \pm 0.9	8.5 \pm 0.3	50.0 \pm 1.7	16.8 \pm 0.4	53.7 \pm 1.5	12.8 \pm 0.5	56.5 \pm 1.2	8.4 \pm 0.4
+ ext	66.3 \pm 1.2	9.8 \pm 0.3	68.1 \pm 1.1	7.7 \pm 0.3	68.9 \pm 1.0	5.5 \pm 0.4	55.2 \pm 1.8	9.6 \pm 0.4	57.7 \pm 1.6	7.4 \pm 0.3	58.7 \pm 1.2	5.4 \pm 0.3



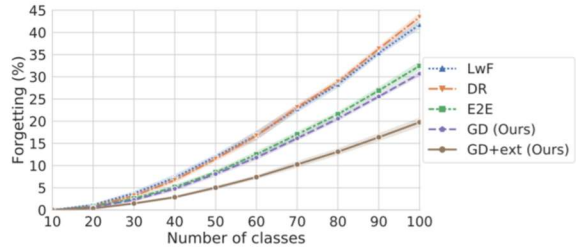
(a) ACC improvement by learning with external data



(b) FGT improvement by learning with external data



(c) ACC with respect to the number of trained classes



(d) FGT with respect to the number of trained classes

Figure 2. Experimental results on CIFAR-100. (a,b) Arrows show the performance gain in the average incremental accuracy (ACC) and average forgetting (FGT) by learning with unlabeled data, respectively. (c,d) Curves show ACC and FGT with respect to the number of trained classes when the task size is 10. We report the average performance of ten trials.

- Effect of the reference models

Table 2. Comparison of models learned with different reference models on CIFAR-100 when the task size is 10. “ \mathcal{P} ,” “ \mathcal{C} ,” and “ \mathcal{Q} ” stand for the previous model, the teacher for the current task, and their ensemble model, respectively.

\mathcal{P}	\mathcal{C}	\mathcal{Q}	ACC (\uparrow)	FGT (\downarrow)
✓			62.9 ± 1.2	14.7 ± 0.4
✓	✓		67.0 ± 0.9	10.7 ± 0.3
		✓	65.7 ± 0.9	11.2 ± 0.2
✓	✓	✓	68.1 ± 1.1	7.7 ± 0.3

- Effect of the teacher for the current task

Table 3. Comparison of models learned with a different teacher for the current task \mathcal{C} on CIFAR-100 when the task size is 10. For “cls,” \mathcal{C} is not trained but the model learns by optimizing the learning objective of \mathcal{C} directly. The model learns with the proposed 3-step learning for “dst.” The confidence loss is added to the learning objective for \mathcal{C} for “cnf.” We do not utilize \mathcal{Q} for this experiment, because “cls” has no explicit \mathcal{C} .

\mathcal{C}	Confidence	ACC (\uparrow)	FGT (\downarrow)
X		62.9 ± 1.2	14.7 ± 0.4
cls		62.9 ± 1.3	14.5 ± 0.5
cls	cnf	65.3 ± 1.0	11.7 ± 0.3
dst		66.2 ± 1.0	11.2 ± 0.3
dst	cnf	67.0 ± 0.9	10.7 ± 0.3

- Effect of balanced fine-tuning

Table 4. Comparison of different balanced learning strategies on CIFAR-100 when the task size is 10. “DW,” “FT-DSet,” and “FT-DW” stand for training with data weighting in Eq. (10) for the entire training, fine-tuning with a training dataset balanced by removing data of the current task, and fine-tuning with data weighting, respectively.

Balancing	ACC (\uparrow)	FGT (\downarrow)
\times	67.1 \pm 0.9	11.5 \pm 0.3
DW	67.9 \pm 0.9	9.6 \pm 0.2
FT-DSet	67.2 \pm 1.1	8.4 \pm 0.2
FT-DW	68.1 \pm 1.1	7.7 \pm 0.3

- Effect of external data sampling

Table 5. Comparison of different external data sampling strategies on CIFAR-100 when the task size is 10. “Prev” and “OOD” columns describe the sampling method for data of previous tasks and out-of-distribution data, where “Pred” and “Random” stand for sampling based on the prediction of the previous model \mathcal{P} and random sampling, respectively. In particular, for when sampling OOD by “Pred,” we sample data minimizing the confidence loss \mathcal{L}_{cnf} . When only Prev or OOD is sampled, the number of sampled data is matched for fair comparison.

Prev	OOD	ACC (\uparrow)	FGT (\downarrow)
\times	\times	65.0 \pm 1.1	12.1 \pm 0.3
\times	Random	67.6 \pm 0.9	9.0 \pm 0.3
Pred	\times	66.0 \pm 1.2	7.8 \pm 0.3
Pred	Pred	65.7 \pm 1.1	10.2 \pm 0.2
Pred	Random	68.1 \pm 1.1	7.7 \pm 0.3

8. Conclusion

- Novel class-incremental learning scheme that uses large stream of unlabeled data
- Global knowledge distillation
- Learning strategy to avoid overfitting to most recent task
- Confidence based sampling method to effectively leverage unlabeled dataset

9. Quiz questions:

9.1 Which of the following statements are true about the global distillation model

- A) Training a reference teacher's model to specialize in learning only the current task
- B) Knowledge distillation for the ensemble model is performed over both the training data and sampled external unlabeled data
- C) Fine-tuning using data weighting is performed over all model parameters
- D) Global distillation model is trained through knowledge distillation over 3 reference models.

Answer: A and D

9.2 Which of the following statements are true about confidence calibration for sampling:

- A) Confidence calibration is performed on all reference models
- B) It prevents the model from making overconfident predictions on OOD data by optimizing over the confidence loss
- C) Confidence calibrated outputs are produced by optimizing the loss function over only the sampled external dataset.
- D) Confidence calibrations increase the overall accuracy of the model by sampling better external data from a stream of unlabeled data

Answer: B and D

9.3 Which external data sampling strategy provides the highest model accuracy:

- A) Random sampling of OOD data and sampling based on predictions of previous model
- B) Only random sampling of OOD data
- C) Sampling based on predictions of previous model and sampling OOD data based on predictions of previous model
- D) No external data sampling.

Answer: A

10: FAQ

Q: About quiz question 2 A, isn't confidence calibration done for all reference models?

A: It is done only on the current model, since at the next stage the current model becomes a part of the previous models, and we don't need to calibrate them again. Only calibration for the current model is enough.