

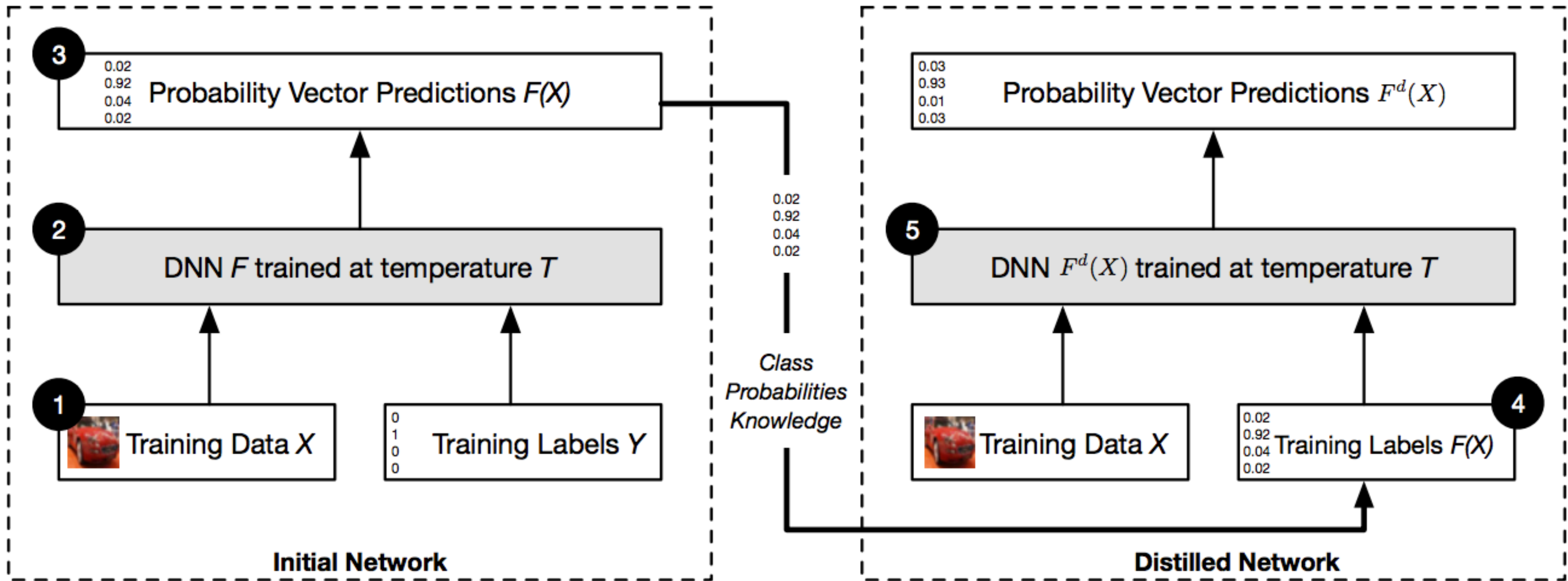




ADVANCE ALL MARCH EVERLASTING

DEFENSE: DISTILLATION

36



An overview of distillation defence mechanism based on a transfer of knowledge contained in probability vectors through distillation: We first train an initial network  $F$  on data  $X$  with a softmax temperature of  $T$ . We then use the probability vector  $F(X)$ , which includes additional knowledge about classes compared to a class label, predicted by network  $F$  to train a distilled network  $F_d$  at temperature  $T$  on the same data  $X$ .

Paper N. McDaniel, P. Mu, X. Jha, S. & Swami, A. (2016, May) Distillation as adversarial perturbation against deep neural networks



## DEFENSE: EXTENDED DEFENSIVE DISTILLATION

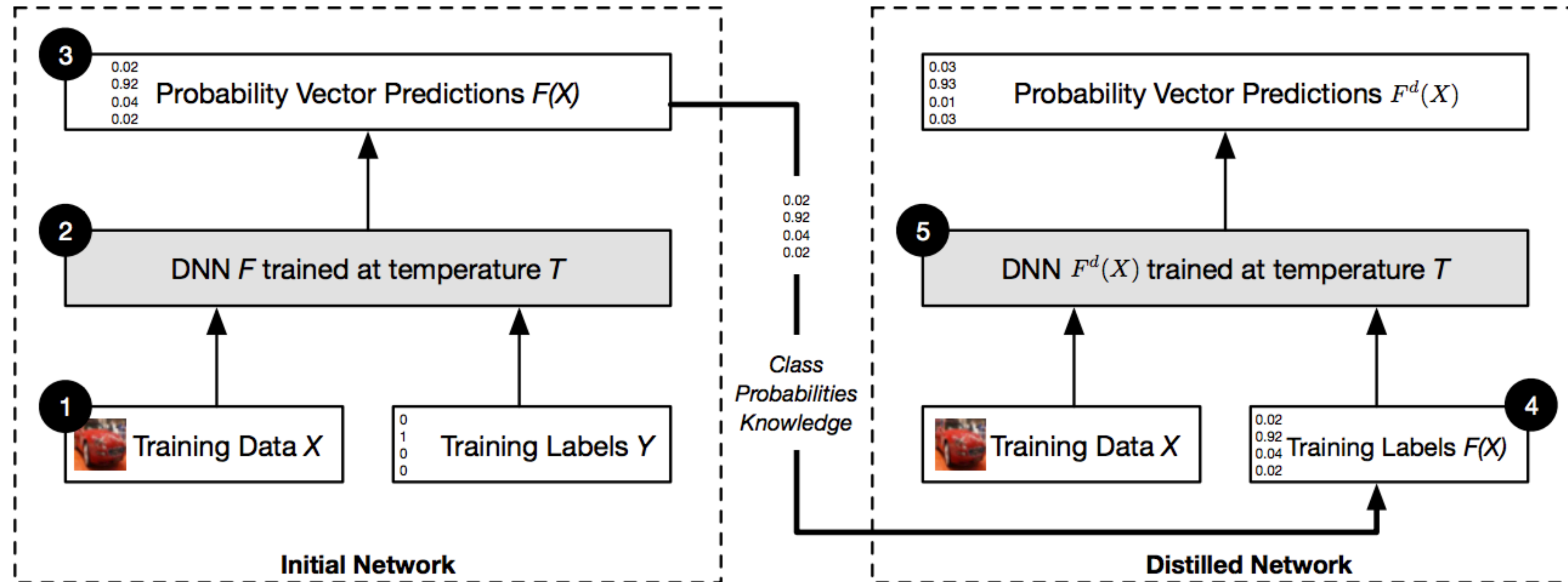
► Uncertainty Measure:

$$\sigma(x) = \frac{1}{N} \sum_{m \in 0..N-1} \left( \sum_{j \in 0...n-1} (z_j^m(x) - \bar{z}_j)^2 \right)$$

► Labelling Vector:

$$k_j(x) = \begin{cases} 1 - \alpha \cdot \frac{\sigma(x)}{\max_{x \in \mathcal{X}} \sigma(x)} & \text{if } j = l \text{ (correct label)} \\ \alpha \cdot \frac{\sigma(x)}{\max_{x \in \mathcal{X}} \sigma(x)} & \text{if } j = n \text{ (outlier class)} \\ 0 & \end{cases}$$

# DEFENSE: DISTILLATION



An overview of distillation defence mechanism based on a transfer of knowledge contained in probability vectors through distillation: We first train an initial network  $F$  on data  $X$  with a softmax temperature of  $T$ . We then use the probability vector  $F(X)$ , which includes additional knowledge about classes compared to a class label, predicted by network  $F$  to train a distilled network  $F^d$  at temperature  $T$  on the same data  $X$ .