

# DINO: distillation with no labels

① DINO: Distillation with No Labels.

Approach: SSL with LB:

Student Network  $g_{\theta_s}$  ; Teacher Net  $g_{\theta_t}$   $\rightarrow$   $k$ -dimensional Output.

softmax conversion of student Output.

$$p_s(x)^i = \frac{\exp \{ g_{\theta_s}(x)^i / \tau_s \}}{\sum_{j=1}^k \exp \{ g_{\theta_s}(x)^j / \tau_s \}} \quad // \text{student}$$

$$p_t(x)^i = \frac{\exp \{ g_{\theta_t}(x)^i / \tau_t \}}{\sum_{j=1}^k \exp \{ g_{\theta_t}(x)^j / \tau_t \}} \quad // \text{teacher.}$$

Optimization goal:  $\min_{\theta_s} H(p_t(x), p_s(x)) \quad // \quad H(a, b) = -a \log b$

Actually: Local-to-Global view

$$\min_{\theta_s} \sum_{x \in \{x_1^g, x_2^g\}} \sum_{\substack{x' \in V \\ x' \neq x}} H(p_t(x), p_s(x'))$$

$V \rightarrow$  set of images with different view of  $x$ .

$$\{ x_1^g, x_2^g, x_1^L, \dots, x_n^L \}$$

$\downarrow$   
global views

Local views  $\{$  smaller denoting  $\}$   
 $\{$  student network  $\}$

$\{$  passed through teacher only  $\}$

(11)

Teacher Network Update:

$$\theta_i^o \leftarrow \lambda \theta_i + (1 - \lambda) \theta_s \quad // \text{Polyak-Ruppert Avg}$$

Interesting update.

Avoid collapse:

centering:  $g_t(x) \leftarrow g_t(x) + c$

$$c \leftarrow mc + (1 - m) \frac{1}{B} \sum_{i=1}^B g_{\theta_t}^{(i)}(x)$$

Exp: Avg