

FIGURE 1. Contrasting supervised, unsupervised and self-supervised learning paradigms for training a model f using raw data x , labels y , and loss function \mathcal{L} . Self-supervision methods introduce pretext tasks \mathcal{P} that generate pseudolabels z for discriminative training of f .

Module – 2

Basics of Foundation Models, SSL formalisms, definitions and examples of

- **Pretext tasks**
- **Losses**
- **Downstream adaptations**
in a domain-agnostic setting

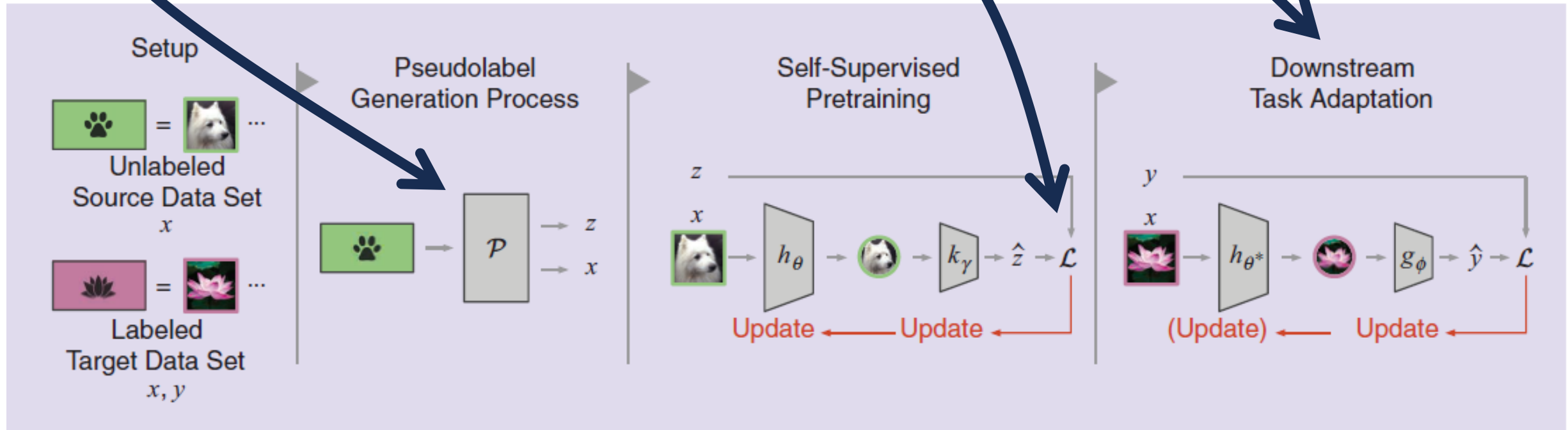


FIGURE 2. The self-supervised workflow starts with an unlabeled source data set and a labeled target data set. As defined by the pretext task, pseudolabels are programmatically generated from the unlabeled set. The resulting inputs, x and pseudolabels z , are used to pretrain the model $k_\gamma(h_\theta(\cdot))$ —composed of feature extractor h_θ and output k_γ modules—to solve the pretext task. After pretraining is complete, the learned weights θ^* of the feature extractor h_{θ^*} are transferred and used together with a new output module g_ϕ to solve the downstream target task.

4.10 Theory

Authors: Aditi Raghunathan, Sang Michael Xie, Ananya Kumar, Niladri Chatterji, Rohan Taori, Tatsunori Hashimoto, Tengyu Ma

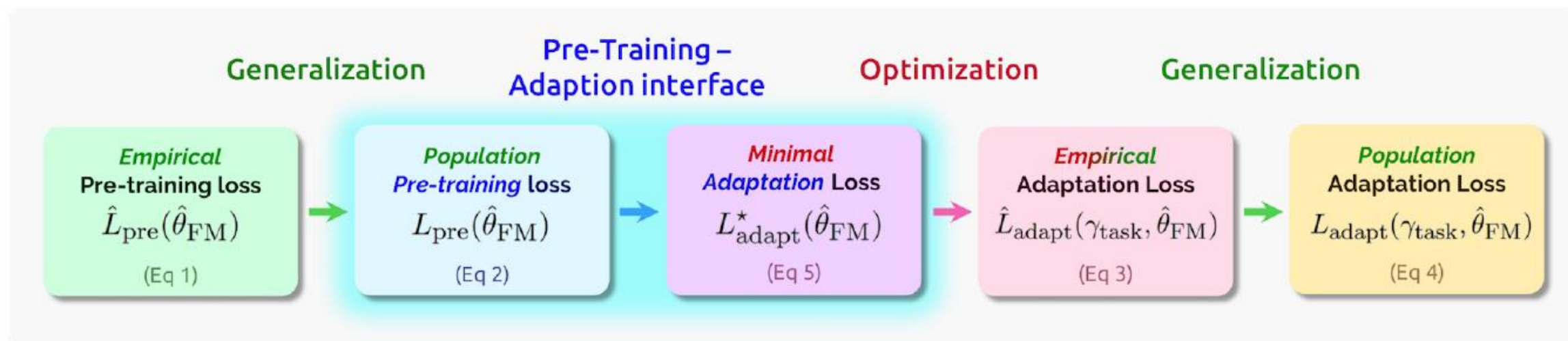
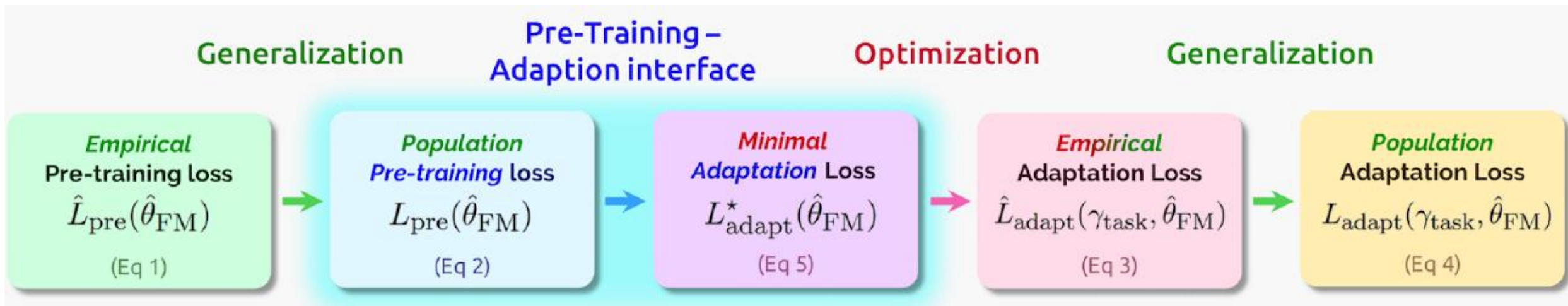


Fig. 22. The analysis of foundation models from pretraining on diverse data to downstream performance on adapted tasks involves capturing the relation between different loss terms as shown above. The main challenge is to analyze the highlighted pretraining-adaptation interface which requires reasoning carefully about the population losses in addition to the model architecture, losses and data distributions of the pretraining and adaptation stages (§4.10.2: THEORY-INTERFACE). Analysis of generalization and optimization largely reduces to their analysis in standard supervised learning.



As shown in Figure 22, the main missing link beyond standard supervised theory is:

Under what conditions does a small population pretraining loss $L_{\text{pre}}(\hat{\theta}_{\text{FM}})$ imply a small minimal adaptation loss $L_{\text{adapt}}^(\hat{\theta}_{\text{FM}})$ and why?*

Supervised Learning.

Target task: t

Supervised (Labeled) Data

$$D_t = \{x_i^{(t)}, y_i^{(t)}\}_{i=1}^N$$

N sample training data.

For each x ($x_i^{(t)}$) $\in D_t$

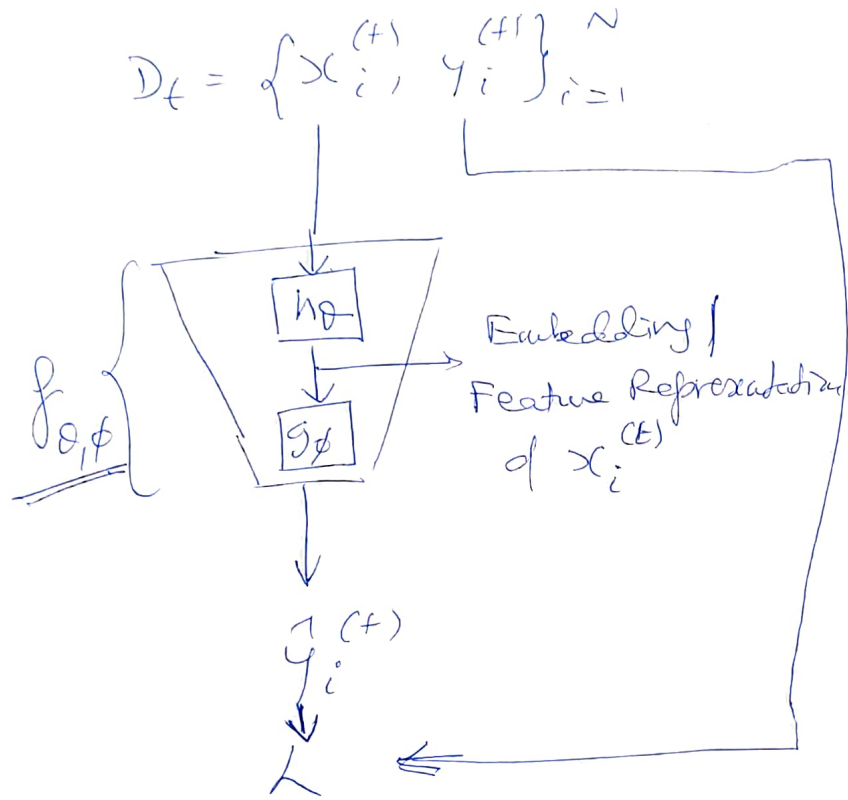
predicted label: $f(x) = y$

$$\underline{\underline{f \rightarrow g_\phi(h_\theta)}}$$

h_θ : Representation feature extractor fn.

g_ϕ : classifier/regression fn.

$$\boxed{\hat{y}_i^{(t)} = g_\phi(h_\theta(x_i^{(t)}))}$$



$$L(\hat{y}_i^{(+)}, y_i^{(+)})$$

$$\downarrow$$

$$L(f(x_i^{(+)}), y_i^{(+)})$$

$$L(g_\phi(h_\theta(x_i^{(+)})), y_i^{(+)})$$

$$\Rightarrow L(i, \theta, \phi)$$

Total loss on D_t

[Full Batch Presentation]

$$L(\theta, \phi) = \sum_{i=1}^N L(i, \theta, \phi)$$

$$f^* = \theta^*, \phi^*$$

$$= \underset{\theta, \phi}{\text{argmin}} L(\theta, \phi)$$

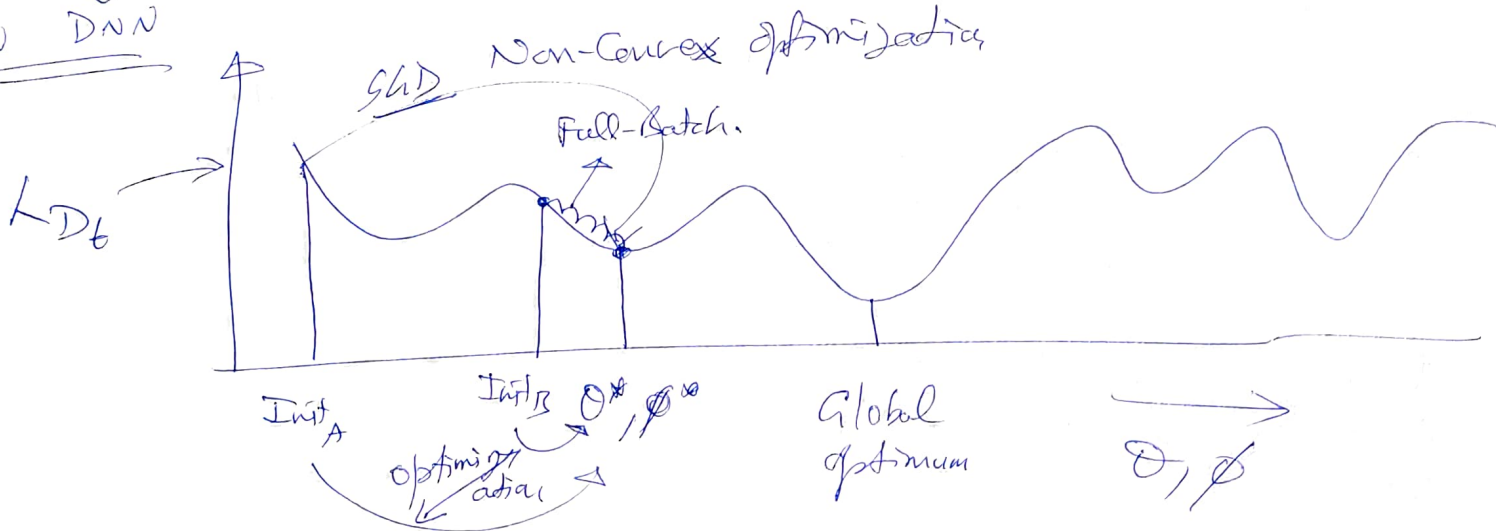
(3)

$$\begin{aligned}
 \theta^* &= \theta^*, \phi^* = \underset{\theta, \phi}{\operatorname{argmin}} L(\theta, \phi) \quad \hat{y}_i^{(+)} \\
 &= \underset{\theta, \phi}{\operatorname{argmin}} \sum_{i=1}^N L\left(\theta_\phi(h_\theta(x_i^{(+)})), y_i^{(+)}\right)
 \end{aligned}$$

$h_\theta, g_\phi \Rightarrow \theta, \phi$: 100-1000s of mtd of parameters
 : Deep backbones

$$L_{D_t} = \sum_{x, y \in D_t} L(\cdot, \cdot)$$

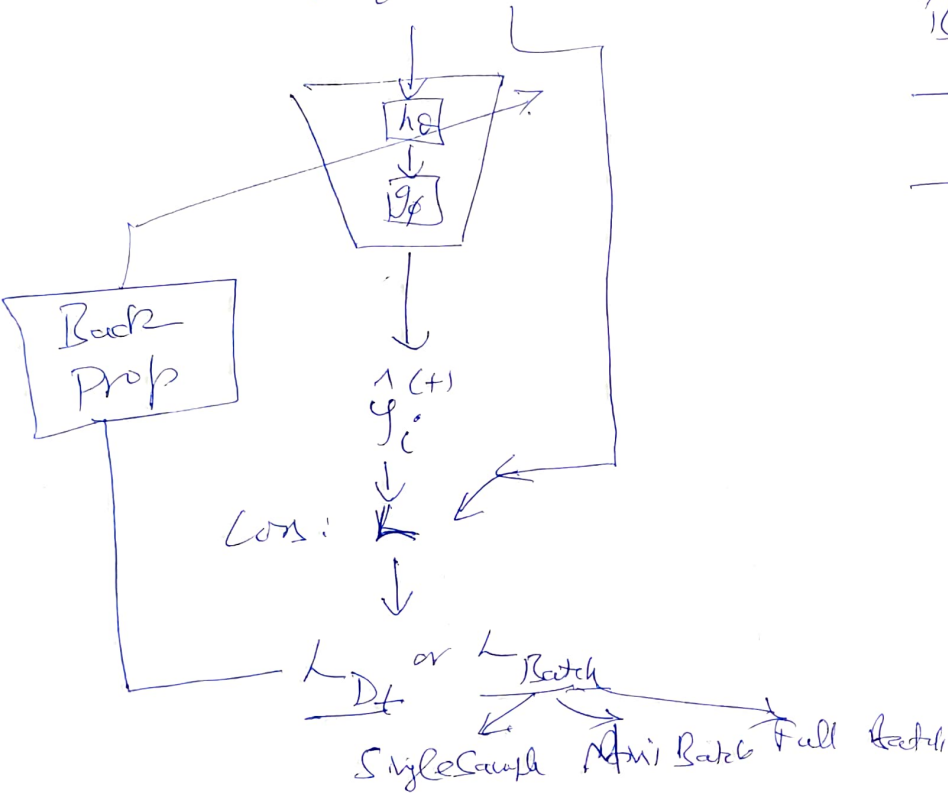
CNN DNN



Finding $\theta^*, \phi^* : f^* \leadsto h_{\theta^*}, g_{\phi^*}$

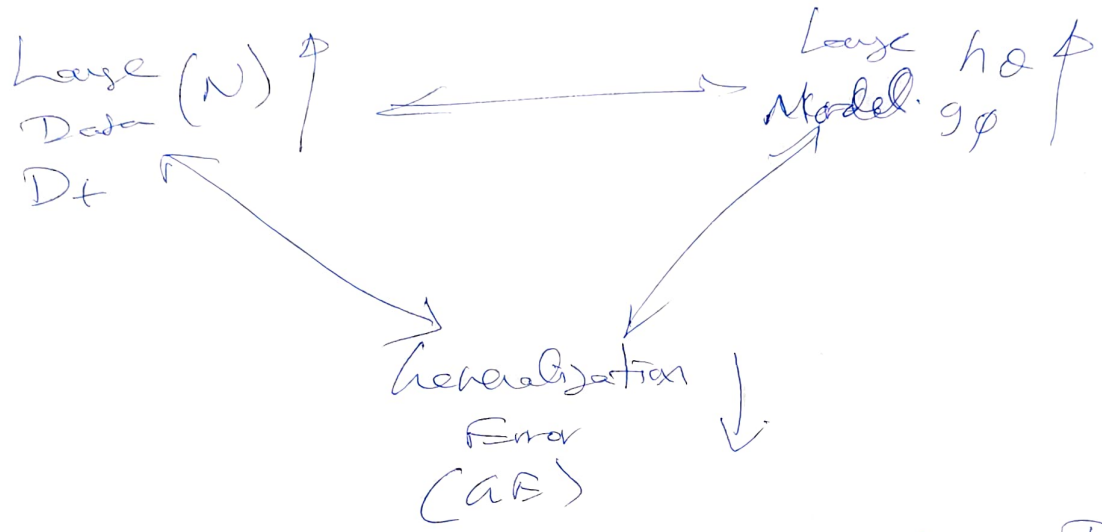
(4)

$$D_t = \{x_i^{(t)}, y_i^{(t)}\}_{i=1}^{N_t}$$



Optimization

By Gradient Descent
 → Back Propagation
 → updated wts using
Ins fn. L_{Batch}



1. Given unseen test/inference \rightarrow Larger D_t [For Low GE]
Larger $D_t \rightarrow$ More Complex Decision Surfaces
 \rightarrow Complex models
2. Given a task \rightarrow VAT \rightarrow Large/Deep Models (More Complex models)
By PAC/Sample Complexity (SC) Bound \rightarrow Larger D_t ($N > SC$)

Unsupervised Learning

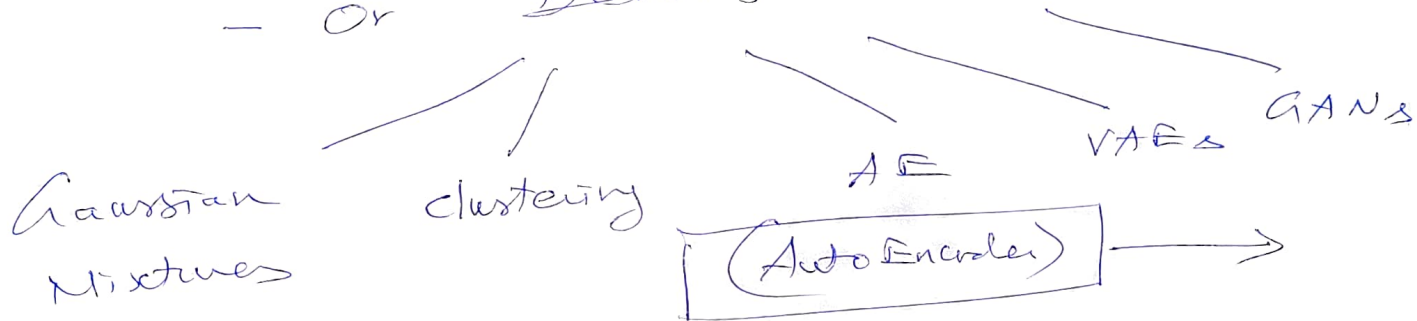
①

□ Learn from 'unlabeled' data

— Learn Compact 'Latent' representations

— Build Generative Models

— Or Density Estimators

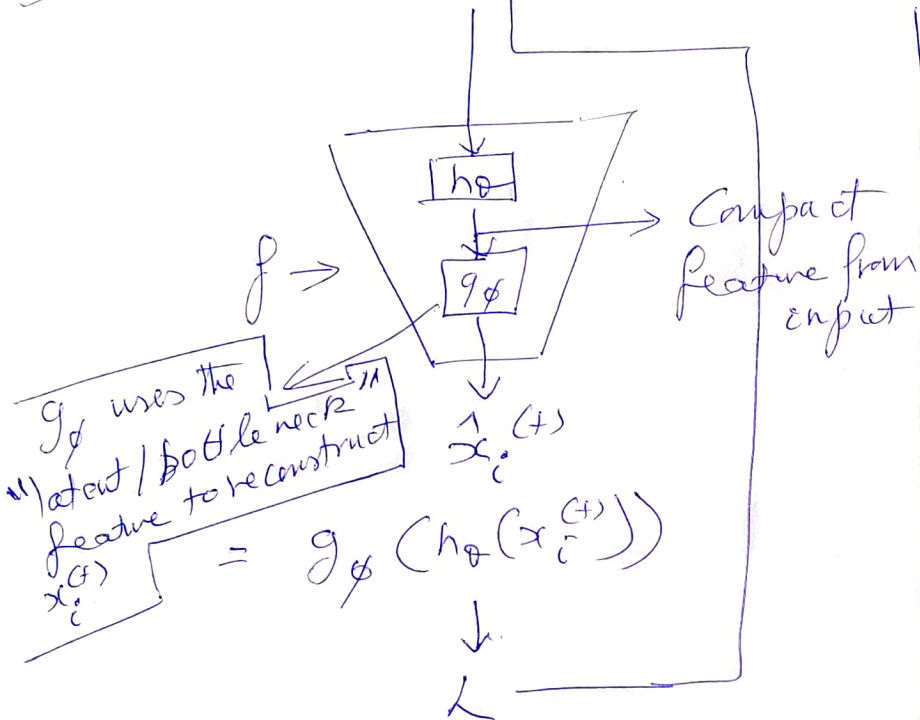


AUTOENCODER

$$D_t = \{x_i^{(+)}\}_{i=1}^N$$

No label - $y_i^{(+)}$

(2)



g_ϕ uses the "latent / bottle neck feature to reconstruct $x_i^{(+)}$ "

$$= g_\phi(h_\theta(x_i^{(+)}))$$

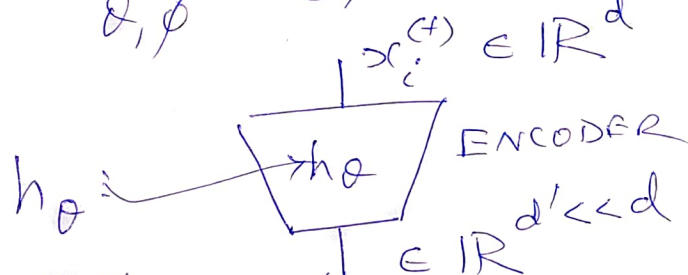
$$L(\hat{x}_i^{(+)}, x_i^{(+)})$$

$$L(g_\phi(h_\theta(x_i^{(+)})), x_i^{(+)})$$

AEs optimize θ, ϕ on "reconstruction" objective

$$\theta^*, \phi^*$$

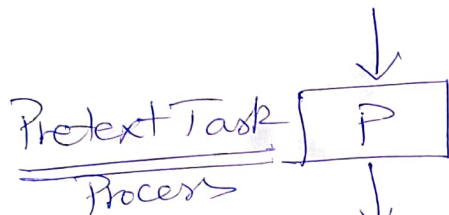
$$= \arg \min_{\theta, \phi} L(g_\phi(h_\theta(x_i^{(+)})), x_i^{(+)})$$



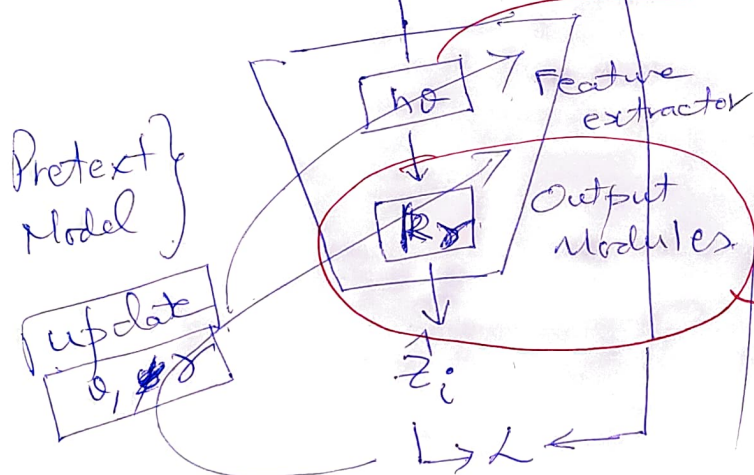
Latent representation

PRETRAINING Self-Supervised Learning DOWN STREAM ①

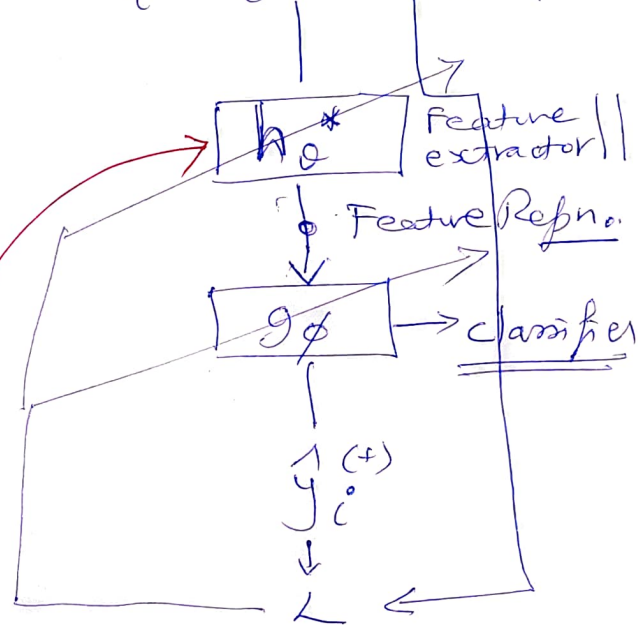
Unlabeled
 $D_S = \{x_i^{(S)}\}_{i=1}^M \quad M \gg N$



$$\bar{D}_S = P(D_S) = \{x_i, z_i\}_{i=1}^M$$



Labeled
 $D_t = \{x_i^{(+)}, y_i^{(+)}\}_{i=1}^N$



Discarded

Linear Readout / Linear Probing

(3)

(ϕ, γ) : Pretrained model weights / Parameters

h_ϕ — Feature extractor

R_γ — Pretext task "specific" head.

Downstream task!

→ Reuse h_ϕ for new task

→ Replace R_γ with a new head g_ϕ designed for the new task.

h_ϕ : Frozen

g_ϕ : Trained/Optimized →

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(g_\phi(h_\phi(x_i^{(t)})), y_i^{(t)})$$

Fine-tuning

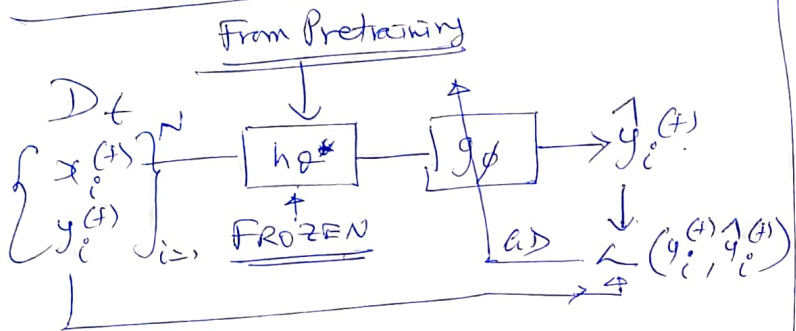
- Retrain the entire network

$$\underline{h_{\theta^*} + g_{\phi}}$$

$$\theta^{**}, \phi^* = \underset{\theta, \phi}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(g_{\phi}(h_{\theta}(x_i^{(t)})), y_i^{(t)})$$

Initialize with θ^*

Linear Readout Probing



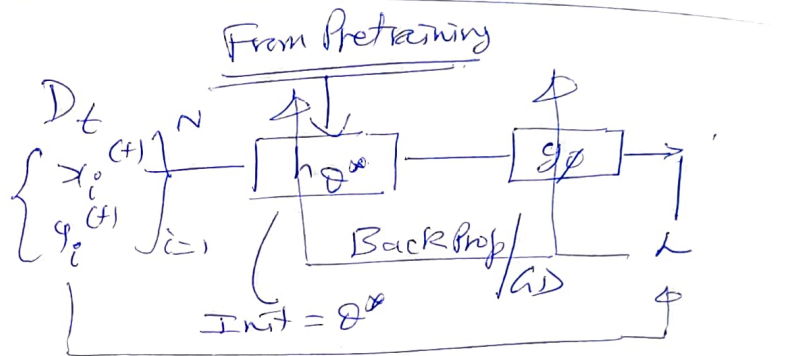
$N \Rightarrow$ Typically "small" \Rightarrow

$N \Rightarrow$ Very Small } Sparse Data regime

- # unique parameters to learn for the target task must be aggressively limited to avoid overfitting.
- g_{ϕ} : Simple Linear In // 1-layer FC

Finetuning

(5)



D_t : Labeled Dataset of Target Task

$N \Rightarrow$ "Moderately more data"

- Enough target data to refine all the model parameters.
- OR when Pretext task/data NOT suited to downstream task.

SSL Scenarios

6

- $N \uparrow$: N is large (Dense labels are available) for the target task.

D_t is large

→ Direct "Supervised learning"

- Best approach.
- SSL may not be helpful

- Domain (Pretraining) \neq Domain (Target)
Very different → Indomain

Indomain
SSL
pretraining
But

eg ImageNet → Radar imaging
SSL: Not very effective

Pretraining Unlabeled Data with SSL
→ Highly effective.

SSL Scenarios (Contd.)

(7)

Domain (Pre Training) \approx Domain (Target)

|
Imagenet

↓
every day objects

SSL — Highly effective.

General Transfer Learning

Supervised
Pretraining

vs

Unsupervised
Pretraining
(SSL)

↓
Superior

SSL Design Choices

(8)

- h_θ frozen, g_ϕ trained (L.R)
- $[h_\theta^*(init) + g_\phi]$ trained (F.T)
- $[D_S \approx D_T]?$ $\|D_t\|$ Academia Performance Comparisons
- $[h_\theta, k_y]$ split in a
 - Pretraining Backbone
 - h_θ : early layers } encode "task/data" agnostic patterns
 - : good for "domain mismatch" conditions
- Pretraining - stopping Criterion [Validation?]
- Downstream Performance \propto Pretrained Model Complexity