

# Lecture Notes for **Neural Networks and Machine Learning**



Self-Supervision, Multi-task,  
and Multi-Modal Learning

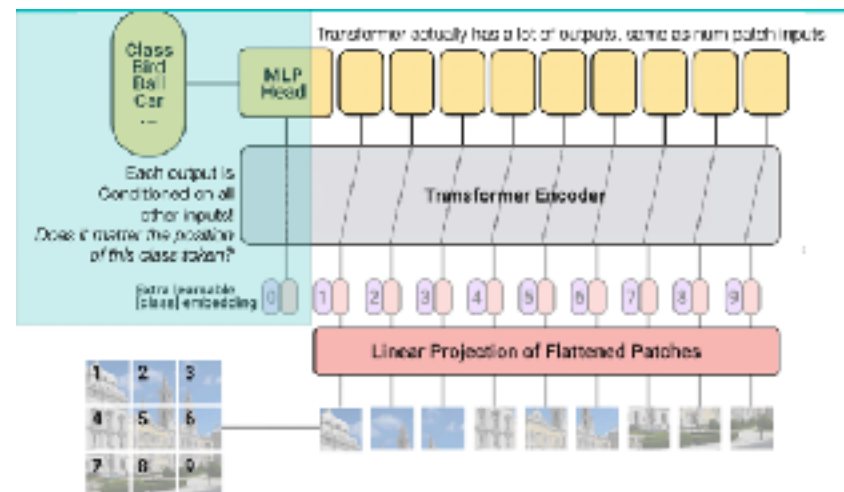


# Logistics and Agenda

- Logistics
  - Grading update
- Agenda
  - Vision transformer and Town Hall
  - Student Paper Presentation
  - Consistency, Contrastive, and Triplet Loss
- Next Time
  - Multi-modal and Multi-Task



# Last Time: Transformers



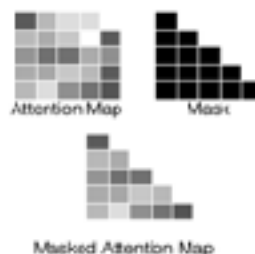
- Decoder only, text encoding happens in attention

$$\mathcal{L}_1(\mathcal{U}) = \sum_{i \in S} \log P(\underbrace{u_i}_{\text{curr}} \mid \underbrace{u_{i-k}, \dots, u_{i-1}}_{\text{other words}}; \underbrace{\mathbf{W}}_{\text{params}})$$

Predict the next word from unlabeled dataset



Also known as  
generative  
pre-training (GPT)



# Vision Transformers Video

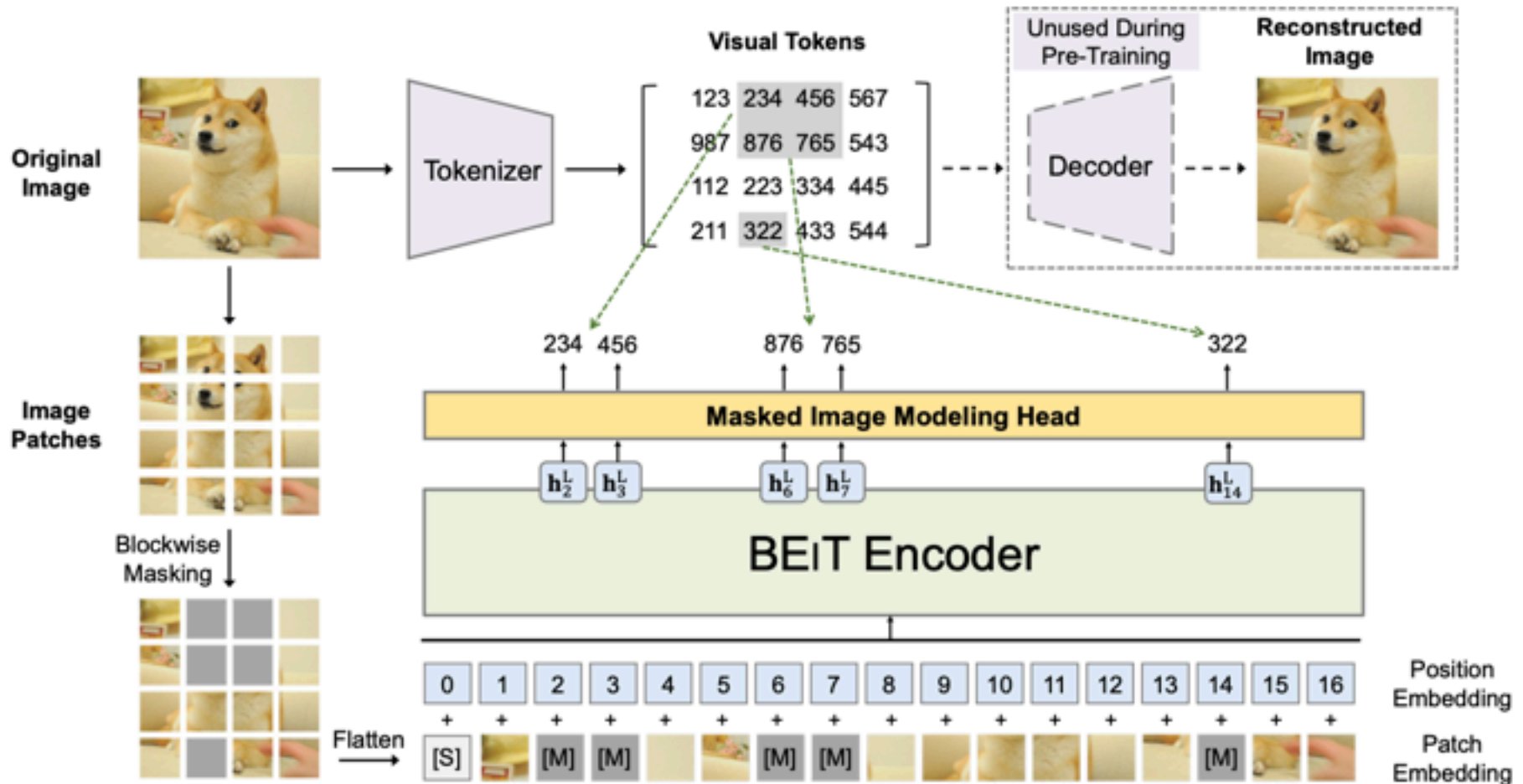


<https://ai.googleblog.com/2020/12/transformers-for-image-recognition-at.html?m=1>

107



# Pre-training for ViT



**Downstream Fine Tuning:** Image classification [CLS], semantic segmentation (patch output), and others...

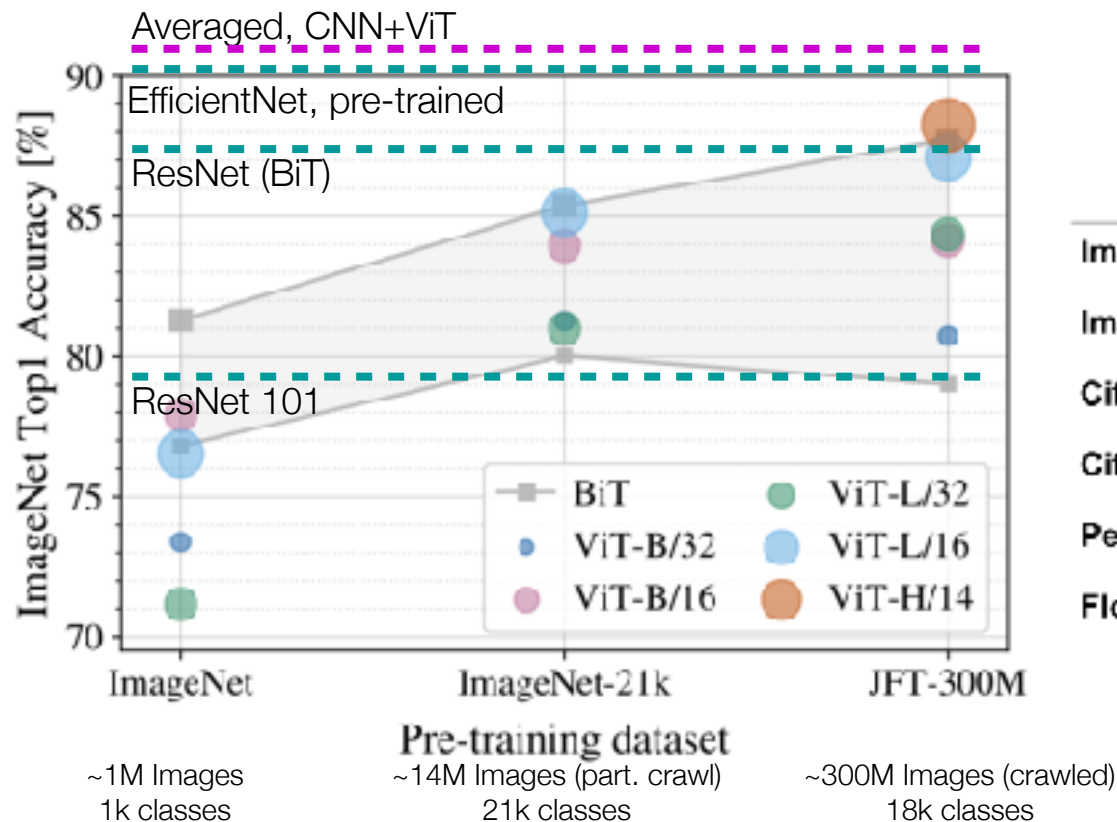
[https://www.fast.ai/2020/01/13/self\\_supervised/](https://www.fast.ai/2020/01/13/self_supervised/)

108



# Fine tuning: Do they work?

- Less than 14M images for pre-training? Do not use as a base model!
  - CNNs will be easier and very performant...



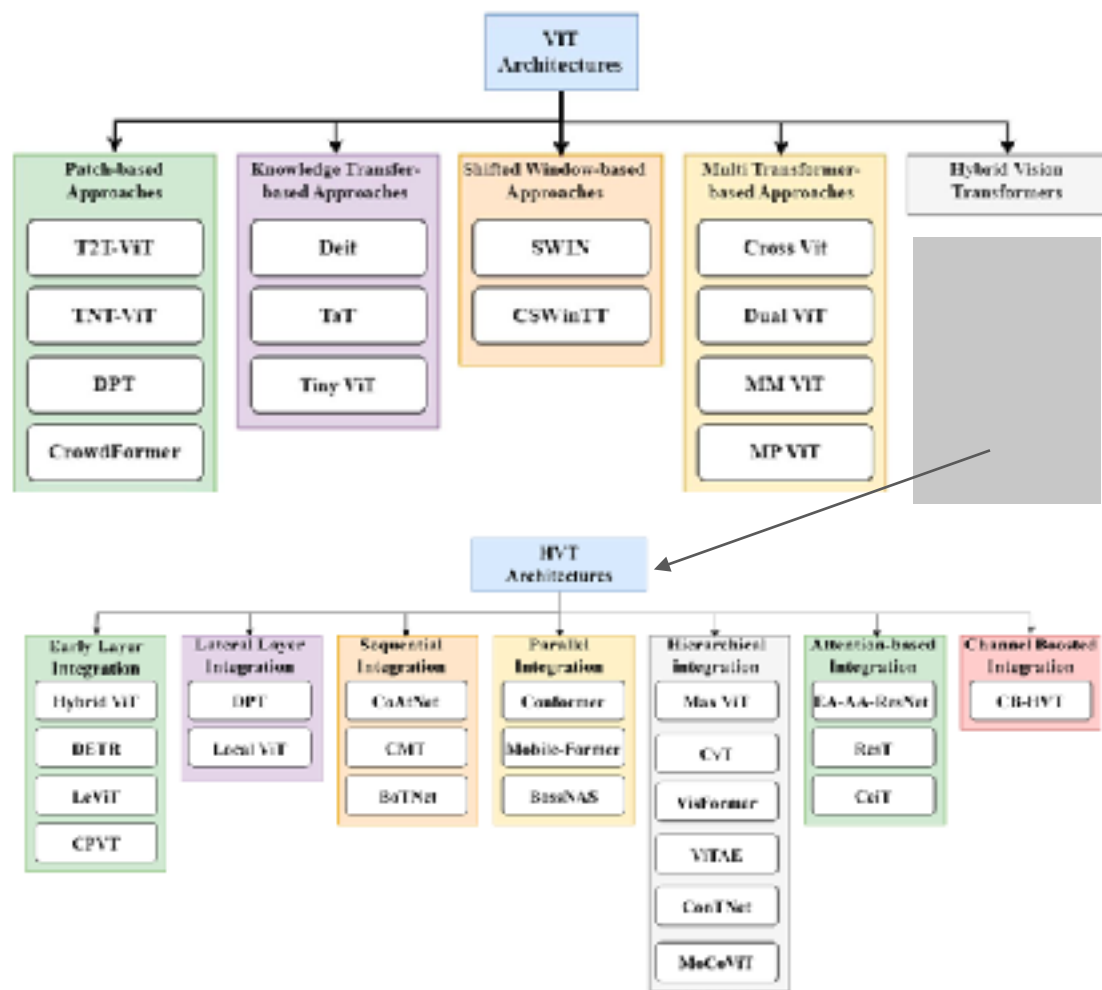
## Transfer Learning From Huge ViT

	BeIT	ViT-H	Previous SOTA
ImageNet	89.5	88.55	← 88.5
ImageNet-Real		90.72	← 90.55
Cifar-10		99.50	← 99.37
Cifar-100	91.8	94.55	← 93.51
Pets		97.56	← 96.62
Flowers		99.68	← 99.63



# Many Variants of the ViT

- ongoing area of research
- Input Patch Structure (overlap)
- Efficient Attention
- Cross Attention (mix of experts)
- Methods of SSL
- Hybrid Approaches
- Image/text generation
  - Generative images patching...



Kahn et al., A survey of the vision transformers and their CNN-transformer variants (2023), *Springer Nature*  
<https://link.springer.com/article/10.1007/s10462-023-10595-0>



# Paper Presentation

## Visual Autoregressive Modeling: Scalable Image Generation via Next-Scale Prediction

Keyu Tian<sup>1,2</sup>, Yi Jiang<sup>2,1</sup>, Zehuan Yuan<sup>2,\*</sup>, Bingyue Peng<sup>2</sup>, Liwei Wang<sup>1,2,\*</sup>

<sup>1</sup>Center for Data Science, Peking University <sup>2</sup>ByteDance Inc.

<sup>2</sup>State Key Lab of General Artificial Intelligence, School of Intelligence Science and Technology, Peking University

keyutian@stu.pku.edu.cn, jiangyi.enjoy@bytedance.com, yuanzehuan@bytedance.com, bingyue.peng@bytedance.com, wanglw@pku.edu.cn

Try and explore our online demo at: <https://var.vision>

Codes and models: <https://github.com/FoundationVision/VAR>



Figure 1: Generated samples from Visual AutoRegressive (VAR) transformers trained on ImageNet. We show 512x512 samples (top), 256x256 samples (middle), and zero-shot image editing results (bottom).





# Transformer and Transfer Learning Town Hall



Hugging Face Text transformers: [https://huggingface.co/transformers/v3.3.1/pretrained\\_models.html](https://huggingface.co/transformers/v3.3.1/pretrained_models.html)

Hugging Face ViT: [https://huggingface.co/docs/transformers/model\\_doc/vit](https://huggingface.co/docs/transformers/model_doc/vit)

Keras text Transformers: [https://keras.io/guides/keras\\_nlp/transformer\\_pretraining/](https://keras.io/guides/keras_nlp/transformer_pretraining/)

Keras ViT: <https://github.com/faustomorales/vit-keras>



# Consistency Loss

I'm from Canada, but live in the States now.

It took me a while to get used to writing boolean variables with an "Is" prefix, instead of the "Eh" suffix that Canadians use when programming.

For example:

```
MyObj.IsVisible
```

```
MyObj.VisibleEh
```



# Unsupervised Consistency Loss

$$\min_{\mathbf{w}} \mathbf{E}_{\mathbf{x}, y \in L} [-\log p_{\mathbf{w}}(y | \mathbf{x})] + \lambda \mathcal{D}_{KL} (p_{\mathbf{w}}(y | \mathbf{x}) || p_{\mathbf{w}}(y | \hat{\mathbf{x}}))$$

no back prop      yes back prop

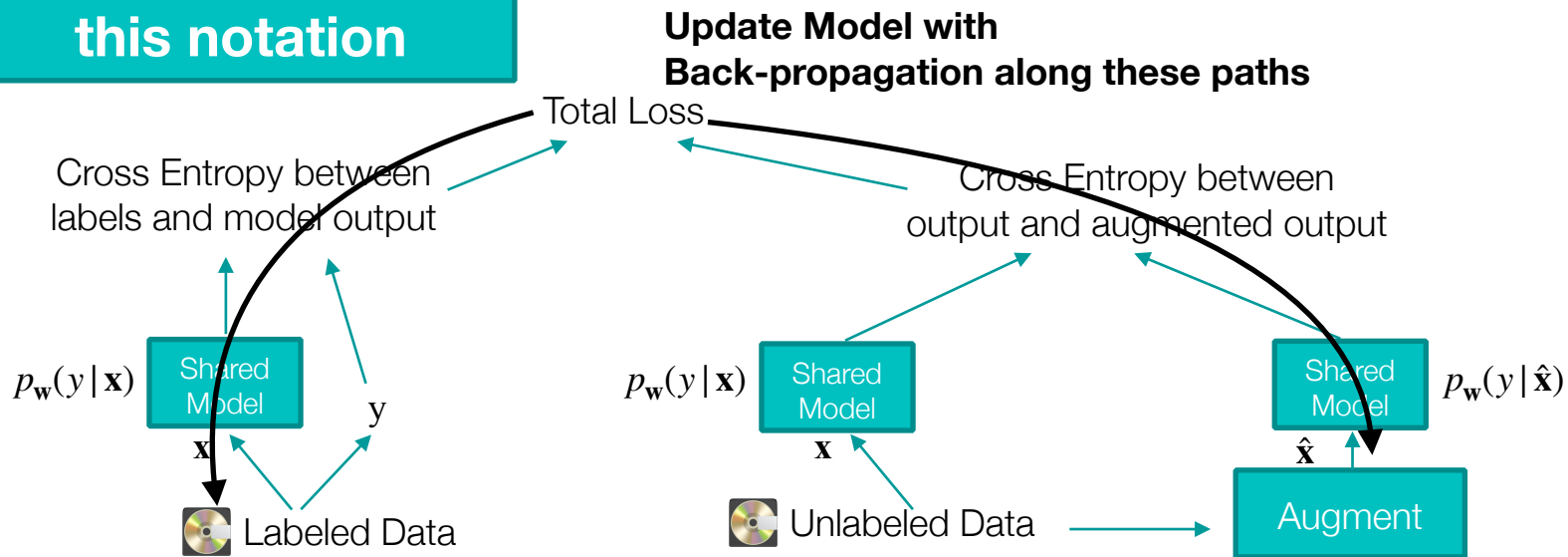
Neural Network approximates  $p(y|\mathbf{x})$  by  $\mathbf{w}$

Sometimes shown as  $p(y | \mathbf{x}; \mathbf{W})$

Use labeled data to minimize network

Sample new  $\mathbf{x}$  from unlabeled pool with function  $q$   
function  $q$  is augmentation procedure  
Minimize cross entropy of two models

**Get accustomed to this notation**



$$\min_{\mathbf{w}} \underbrace{\mathbf{E}_{\mathbf{x}, y \in L} [-\log p_{\mathbf{w}}(y | \mathbf{x})]}_{\text{cross entropy}} + \lambda \underbrace{\mathcal{D}_{KL}(p_{\mathbf{w}}(y | \mathbf{x}) || p_{\mathbf{w}}(y | \hat{\mathbf{x}}))}_{\text{consistency in augmentation}}$$


---

$$E[g] = \sum p(g) \cdot g \quad \text{definition of expected value}$$

$$E[-\log p_{\mathbf{w}}(y | \mathbf{x})] = - \sum p(y) \cdot \log p_{\mathbf{w}}(y | \mathbf{x}) \quad \text{insert -log probability, log likelihood}$$

$$NLL(y, p_{\mathbf{w}}(y | \mathbf{x})) = - \sum_c p(y = c) \cdot \log p_{\mathbf{w}}(y = c | \mathbf{x}) \quad \text{negative log likelihood, discrete classes}$$

$$CE(f, g) = - \sum f(x) \cdot \log g(x) \quad \text{cross entropy of two functions}$$

$$CE(y, p_{\mathbf{w}}(y | \mathbf{x})) = - \sum_c p(y = c) \cdot \log p_{\mathbf{w}}(y = c | \mathbf{x})$$

if  $y=c$  is a probability, these are same  
**NLL is Cross Entropy**

```
cce = tf.keras.losses.CategoricalCrossentropy()
cce(y_true, y_pred)
```



$$\min_{\mathbf{w}} \underbrace{\mathbf{E}_{\mathbf{x}, y \in L} [-\log p_{\mathbf{w}}(y | \mathbf{x})]}_{\text{cross entropy}} + \lambda \underbrace{\mathcal{D}_{KL}(p_{\mathbf{w}}(y | \mathbf{x}) || p_{\mathbf{w}}(y | \hat{\mathbf{x}}))}_{\text{consistency in augmentation}}$$


---

$$\mathcal{D}_{KL}(f || g) = - \sum f(x) \cdot \log \frac{g(x)}{f(x)} \quad \text{definition of Kullback-Leibler (KL) Divergence}$$

$$\mathcal{D}_{KL}(p_{\mathbf{w}}(y | \mathbf{x}) || p_{\mathbf{w}}(y | \hat{\mathbf{x}}))$$

$$\mathcal{D}_{KL}(p(y | \mathbf{x}) || p(y | \hat{\mathbf{x}})) = - \sum p(y | \mathbf{x}) \cdot \log \frac{p(y | \hat{\mathbf{x}})}{p(y | \mathbf{x})} = - \sum p(y | \mathbf{x}) \cdot (\log p(y | \hat{\mathbf{x}}) - \log p(y | \mathbf{x}))$$

$$= - \sum p(y | \mathbf{x}) \cdot \log p(y | \hat{\mathbf{x}}) + \sum p(y | \mathbf{x}) \cdot \log p(y | \mathbf{x})$$

cross entropy of  
augmented and not augmented  
model

global entropy of model  
constant as  $p(y | \mathbf{x})$  has no uncertainty

constant

$$CE(p_{\mathbf{w}}(y | \mathbf{x}), p_{\mathbf{w}}(y | \hat{\mathbf{x}})) = - \sum_c p_{\mathbf{w}}(y = c | \mathbf{x}) \cdot \log p_{\mathbf{w}}(y = c | \hat{\mathbf{x}})$$

```
cce = tf.keras.losses.CategoricalCrossentropy()  
cce(y_pred, y_pred_augmented)
```

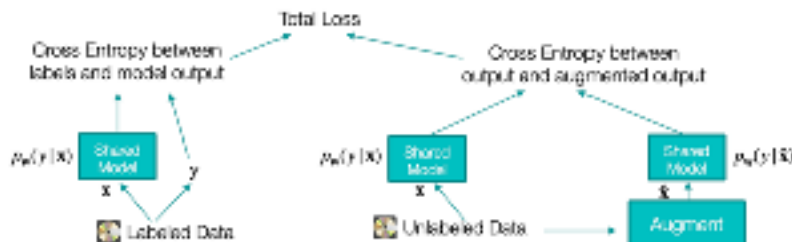


# Aside:

- We have just seen two motivations:

Neural Network approximates  $p(y|x)$  by  $w$   
 Use labeled data to minimize network

Sample new  $x$  from unlabeled pool with function  $g$   
 function  $g$  is augmentation procedure  
 Minimize cross entropy of two models



## intuition of final model

keep labels consistent, any measure would be okay

$$\begin{aligned}
 \mathcal{D}_{KL}(f||g) &= - \sum f(x) \cdot \log \frac{g(x)}{f(x)} \quad \text{definition of Kullback-Leibler (KL) Divergence} \\
 \mathcal{D}_{KL}(p_w(y|x)||p_w(y|\hat{x})) &= - \sum p(y|x) \cdot \log \frac{p(y|\hat{x})}{p(y|x)} = - \sum p(y|x) \cdot (\log p(y|\hat{x}) - \log p(y|x)) \\
 &= - \sum p(y|x) \cdot \log p(y|\hat{x}) + \sum p(y|x) \cdot \log p(y|x) \\
 &\quad \text{cross entropy of augmented and not augmented model} \quad \text{global entropy of model output constant as } p(y|x) \text{ has no uncertainty}
 \end{aligned}$$

```
cce = tf.keras.losses.CategoricalCrossentropy()
cce(y_pred, y_pred_augmented)
```

## mathematics with strict assumptions

not necessarily intuitive, but provides guarantees

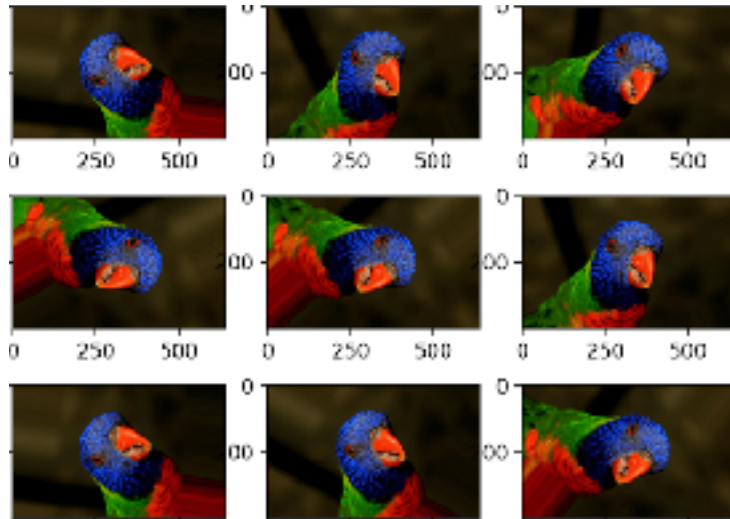
$$\min_w \underbrace{\mathbb{E}_{x,y \in L} [-\log p_w(y|x)]}_{\text{cross entropy}} + \lambda \underbrace{\mathcal{D}_{KL}(p_w(y|x)||p_w(y|\hat{x}))}_{\text{consistency in augmentation}}$$



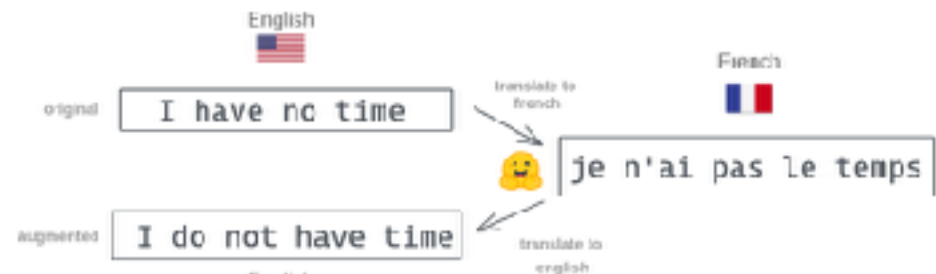
# Augmentation with Consistency Loss

$$\min_{\mathbf{w}} \underbrace{\mathbb{E}_{\mathbf{x}, y \in L} [-\log p_{\mathbf{w}}(y | \mathbf{x})]}_{\text{cross entropy}} + \lambda$$

$$\underbrace{\mathcal{D}_{KL}(p_{\mathbf{w}}(y | \mathbf{x}) || p_{\mathbf{w}}(y | \hat{\mathbf{x}}))}_{\text{consistency in augmentation}}$$



Synonym Replacement



Back Translation



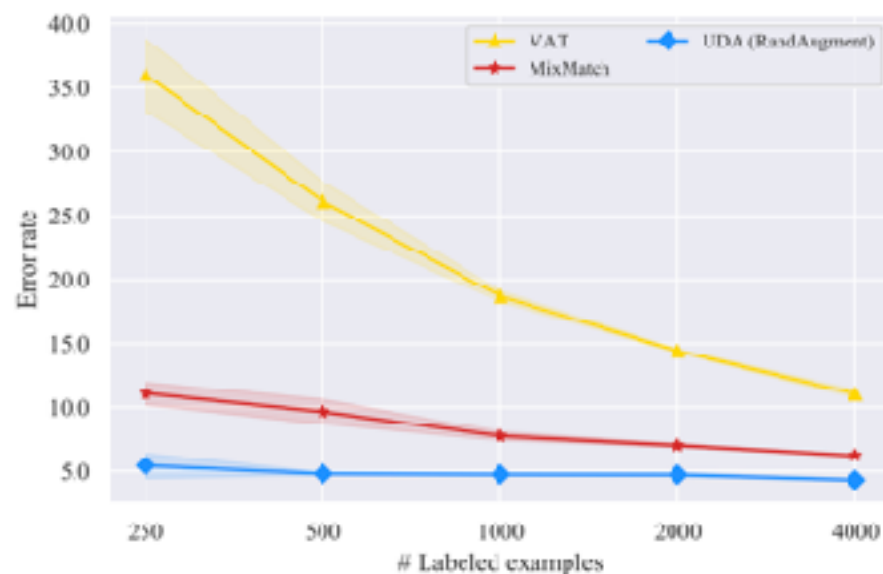
# Unsupervised Consistency Loss

Augmentation (# Sup examples)	Sup (50k)	Semi-Sup (4k)
Crop & flip	5.36	10.94
Cutout	4.42	5.43
RandAugment	<b>4.23</b>	<b>4.32</b>

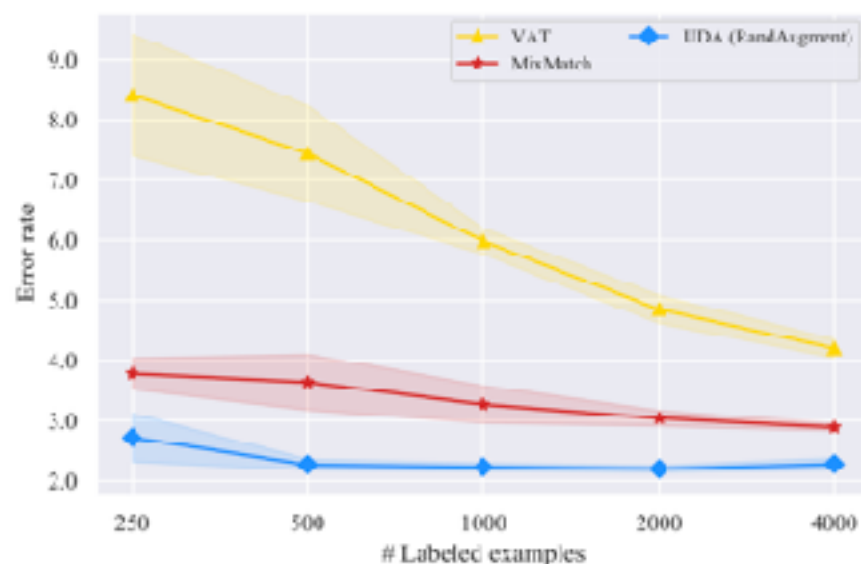
Table 1: Error rates on CIFAR-10.

Augmentation (# Sup examples)	Sup (650k)	Semi-sup (2.5k)
<del>X</del>	38.36	50.80
Switchout	37.24	43.38
Back-translation	<b>36.71</b>	<b>41.35</b>

Table 2: Error rate on Yelp-5.



(a) CIFAR-10



(b) SVHN





# Unsupervised Consistency Loss

Semi-supervised setting

Initialization	UDA	IMDb (20)	Yelp-2 (20)	Yelp-5 (2.5k)	Amazon-2 (20)	Amazon-5 (2.5k)	DBpedia (140)
Random	$\times$ ✓	43.27 25.23	40.25 8.33	50.80 41.35	45.39 16.16	55.70 44.19	41.14 7.24
BERT <sub>BASE</sub>	$\times$ ✓	18.40 5.45	13.60 2.61	41.00 33.80	26.75 3.96	44.09 38.40	2.58 1.33
BERT <sub>LARGE</sub>	$\times$ ✓	11.72 4.78	10.55 2.50	38.90 33.54	15.54 3.93	42.30 37.80	1.68 1.09
BERT <sub>FINETUNE</sub>	$\times$ ✓	6.50 <b>4.20</b>	2.94 <b>2.05</b>	32.39 <b>32.08</b>	12.17 <b>3.50</b>	37.32 <b>37.12</b>	- -



# Contrastive Loss



Jürgen Schmidhuber ✓  
@SchmidhuberAI



DeepSeek [1] uses elements of the 2015 reinforcement learning prompt engineer [2] and its 2018 refinement [3] which collapses the RL machine and world model of [2] into a single net through the neural net distillation procedure of 1991 [4]; a distilled chain of thought system.

REFERENCES (easy to find on the web):

[1] [#DeepSeekR1](#) (2025): Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. arXiv 2501.12948

[2] J. Schmidhuber (JS, 2015). On Learning to Think: Algorithmic Information Theory for Novel Combinations of Reinforcement Learning Controllers and Recurrent Neural World Models. arXiv 1210.0118. Sec. 5.3 describes the reinforcement learning (RL) prompt engineer which learns to actively and iteratively query its model for abstract reasoning and planning and decision making.

[3] JS (2018). One Big Net For Everything. arXiv 1802.08864. See also US11353886B2. This paper collapses the reinforcement learner and the world model of [2] (e.g., a foundation model) into a single network, using the neural network distillation procedure of 1991 [4]. Essentially what's now called an RL "Chain of Thought" system, where subsequent improvements are continually distilled into a single net. See also [5].



# Dealing with Data Sparsity

- How can we augment the number of gradient updates we have, without needing more data collection?
  - Can we leverage the existing examples in the dataset?
  - Can the process provide an exponential number of new examples for back prop?
  - Can it make latent space better behaved for classification?
- **Contrastive Loss:**
  - Use a metric to measure similarity of samples within latent space (e.g., cosine distance, euclidean, etc.)
  - Randomly sample from two or more classes
  - Push same classes together (within a threshold)
  - Push different classes apart (within a threshold)

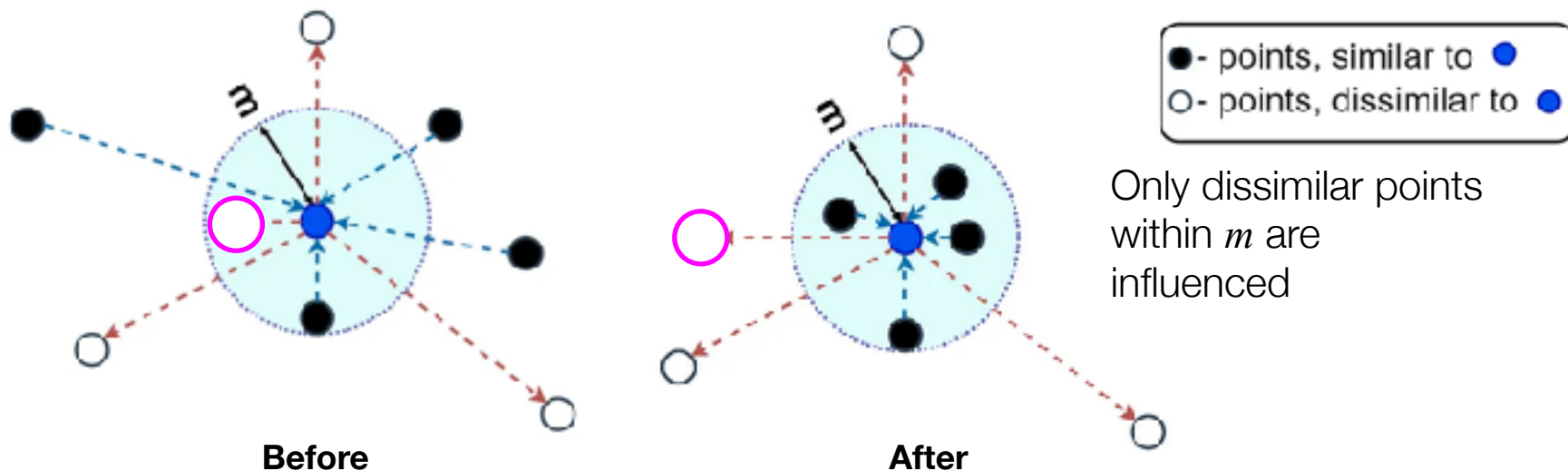


# Contrastive loss

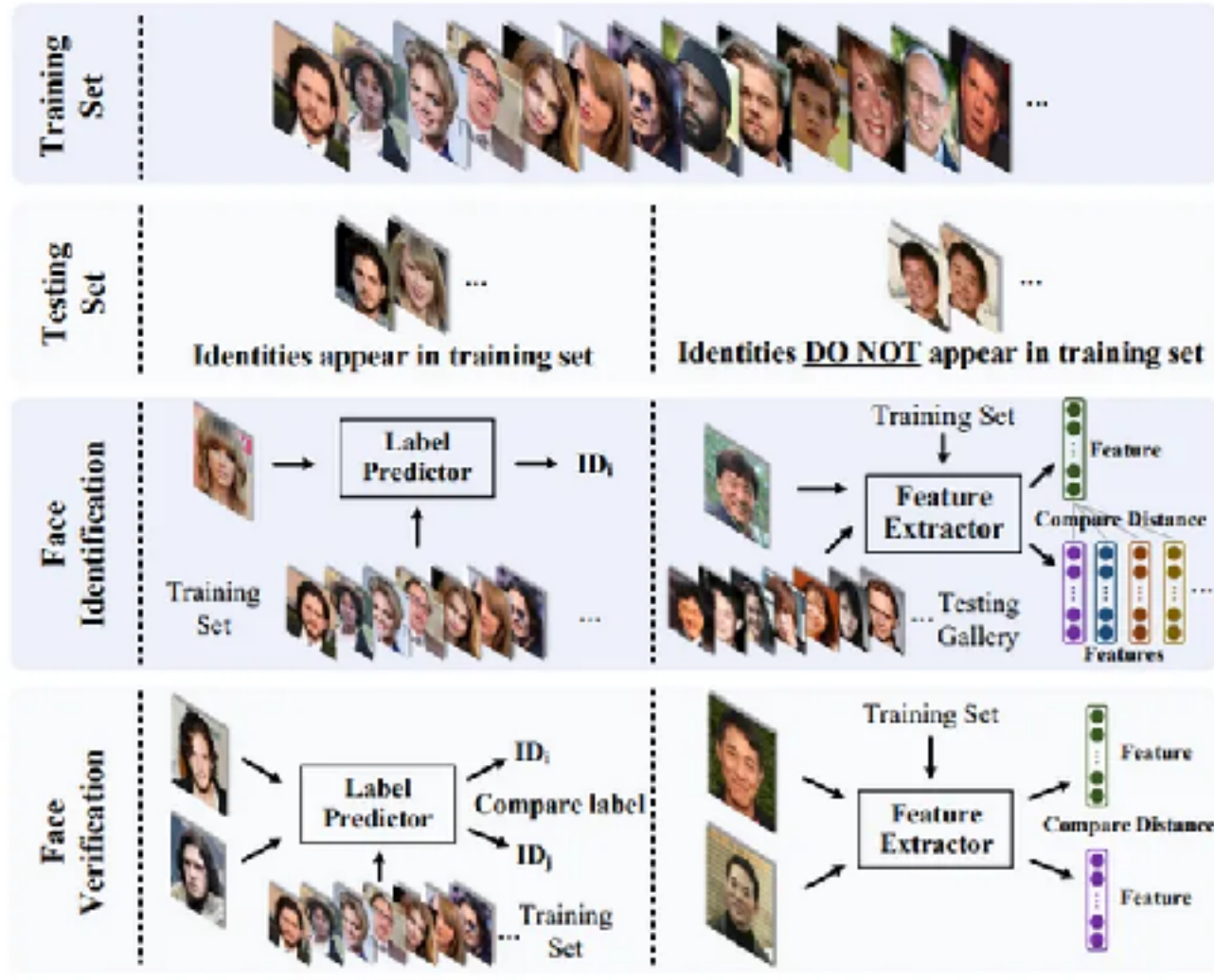
Latent representations of  $\mathbf{x}$ , from model  $f$

$$D(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \|f_{\theta}(\mathbf{x}^{(i)}) - f_{\theta}(\mathbf{x}^{(j)})\|_2$$

$$\mathcal{L}_c = \underbrace{\sum_{i,j \in S} D(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})}_{\text{similar}} + \underbrace{\sum_{i,j \in \hat{S}} \max(0, m - D(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}))}_{\text{not similar}}$$



# Contrastive loss in Face Detect/Identify



Contrastive Loss great for authenticating without re-training

However, the pairs chosen for positive and negative samples tends to be very sensitive to sampling for good performance.

