

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}

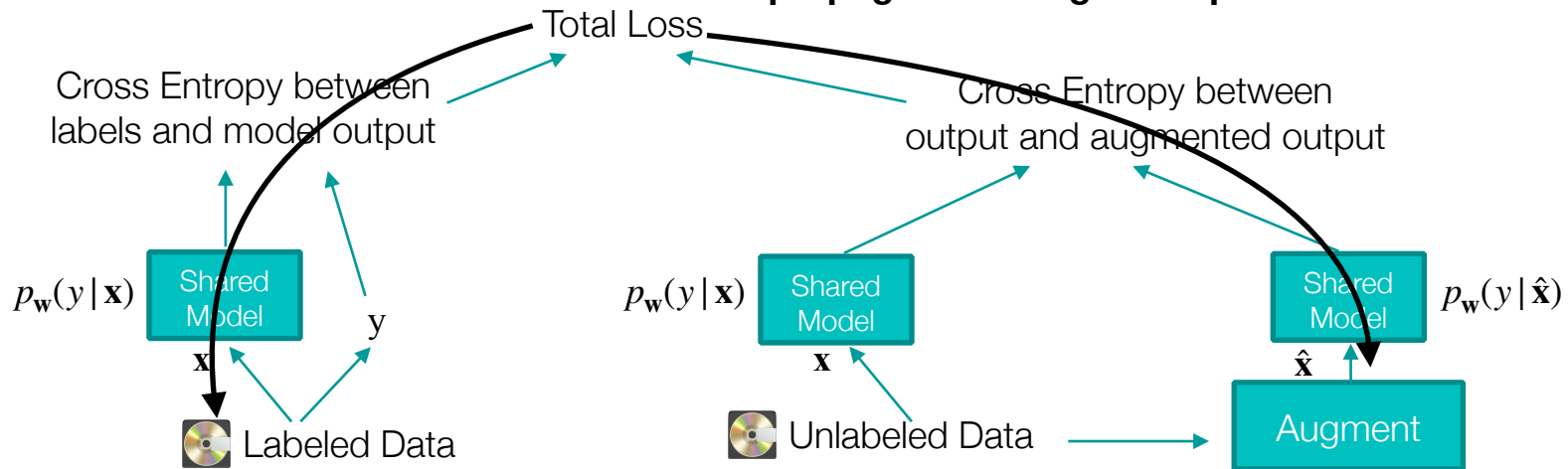
Sometimes shown as $p(y | \mathbf{x}; \mathbf{W})$

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

Update Model with Back-propagation along these paths



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

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

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

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

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

$$CE(y, p_{\mathbf{w}}(y | \mathbf{x})) = - \sum_c (y = c) \cdot \log p_{\mathbf{w}}(y = c | \mathbf{x}) \quad \text{if } y=c \text{ is a probability, these are same equation}$$

```
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

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

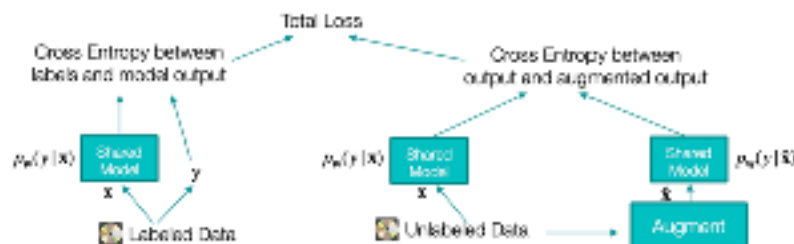


Aside:

- 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 product

keep labels consistent, any measure would be okay

$$\begin{aligned} \mathcal{D}_{KL}(f||g) &= - \sum f(x) \cdot \log \frac{g(x)}{f(x)} \quad \text{definition of Kullback-Leibler (KL) Divergence} \\ \mathcal{D}_{KL}(p_w(y|x)||p_w(y|\hat{x})) &= - \sum p(y|x) \cdot \log \frac{p(y|\hat{x})}{p(y|x)} = - \sum p(y|x) \cdot (\log p(y|\hat{x}) - \log p(y|x)) \\ &= - \sum p(y|x) \cdot \log p(y|\hat{x}) + \sum p(y|x) \cdot \log p(y|x) \end{aligned}$$

cross entropy of augmented and not augmented model global entropy of model output constant as $p(y|x)$ has no uncertainty

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

mathematics with strict assumptions

not necessarily intuitive, but provides guarantees

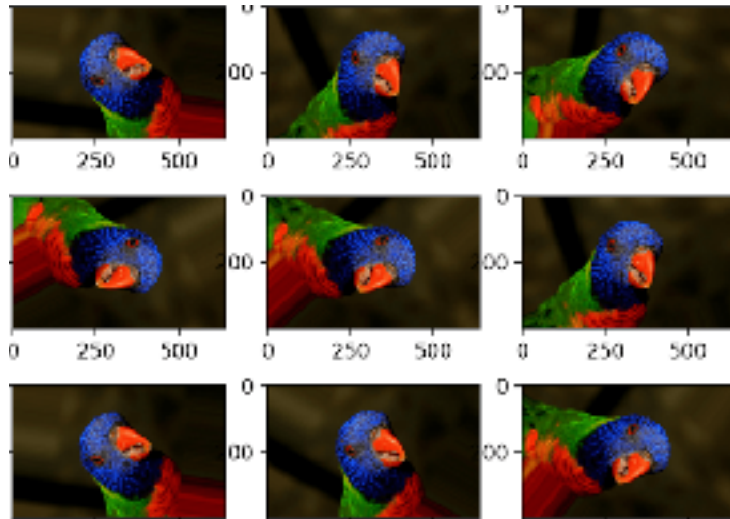
$$\min_w \underbrace{\mathbb{E}_{x,y \in L} [-\log p_w(y|x)]}_{\text{cross entropy}} + \lambda \underbrace{\mathcal{D}_{KL}(p_w(y|x)||p_w(y|\hat{x}))}_{\text{consistency in augmentation}}$$



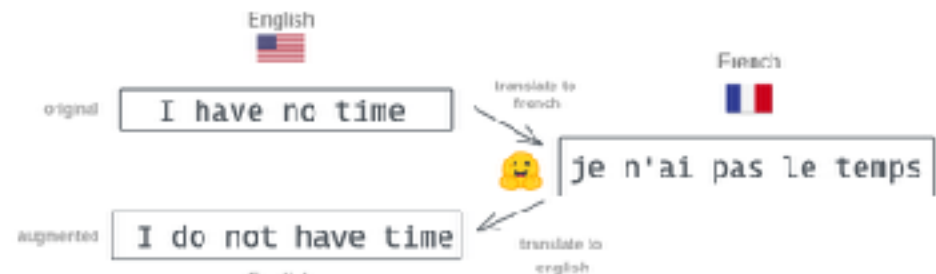
Augmentation with Consistency Loss

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

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



Synonym Replacement



Back Translation



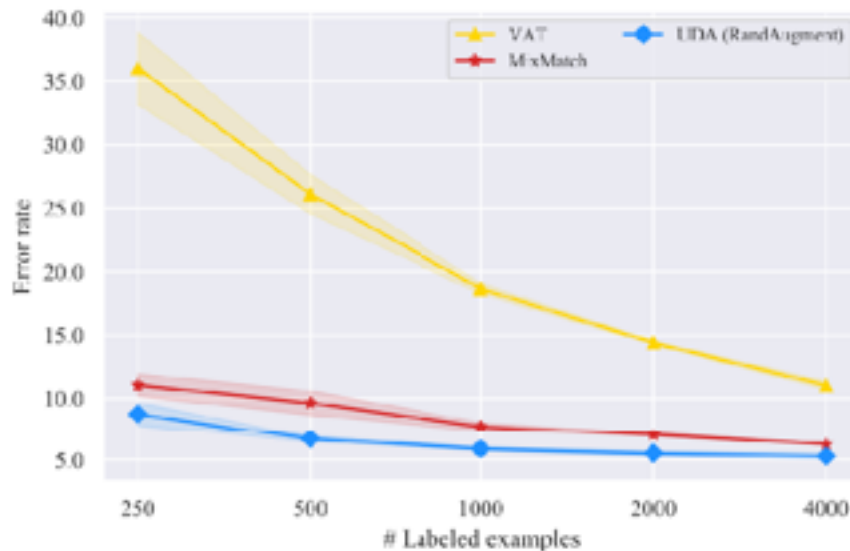
Unsupervised Consistency Loss

Augmentation (# Sup examples)	Sup (50k)	Semi-Sup (4k)
Crop & flip	5.36	16.17
Cutout	4.42	6.42
RandAugment	4.23	5.29

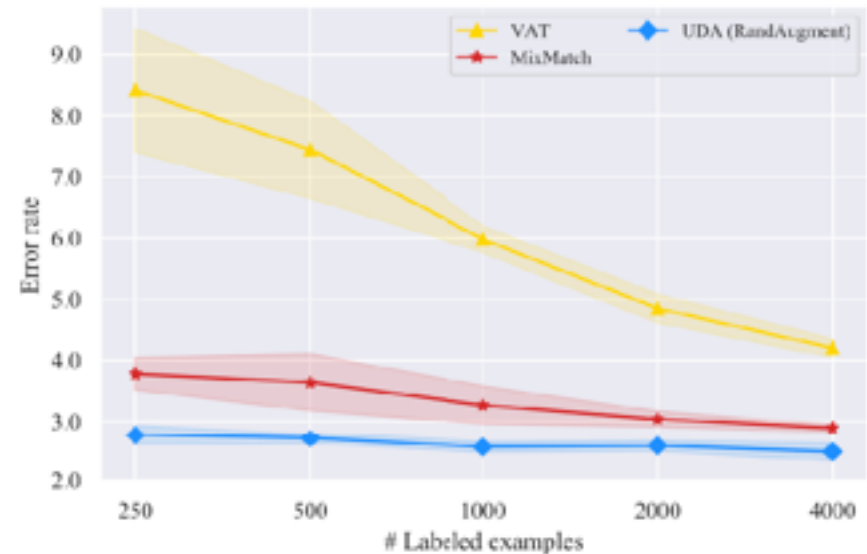
Table 1: Error rates on CIFAR-10.

Augmentation (# Sup examples)	Sup (650k)	Semi-sup (2.5k)
\times	38.36	50.80
Switchout	37.24	43.38
Back-translation	36.71	41.35

Table 2: Error rate on Yelp-5.



(a) CIFAR-10



(b) SVHN



Unsupervised Consistency Loss

Method	Model	# Param	CIFAR-10 (4k)	SVHN (1k)
II-Model (Laine & Aila, 2016)	Conv-Large	3.1M	12.36 ± 0.31	4.82 ± 0.17
Mean Teacher (Tarvainen & Valpola, 2017)	Conv-Large	3.1M	12.31 ± 0.28	3.95 ± 0.19
VAT + EntMin (Miyato et al., 2018)	Conv-Large	3.1M	10.55 ± 0.05	3.86 ± 0.11
SNTG (Luo et al., 2018)	Conv-Large	3.1M	10.93 ± 0.14	3.86 ± 0.27
VAdD (Park et al., 2018)	Conv-Large	3.1M	11.32 ± 0.11	4.16 ± 0.08
Fast-SWA (Athiwaratkun et al., 2018)	Conv-Large	3.1M	9.05	-
ICT (Verma et al., 2019)	Conv-Large	3.1M	7.29 ± 0.02	3.89 ± 0.04
Pseudo-Label (Lee, 2013)	WRN-28-2	1.5M	16.21 ± 0.11	7.62 ± 0.29
LGA + VAT (Jackson & Schulman, 2019)	WRN-28-2	1.5M	12.06 ± 0.19	6.58 ± 0.36
mixmixup (Hataya & Nakayama, 2019)	WRN-28-2	1.5M	10	-
ICT (Verma et al., 2019)	WRN-28-2	1.5M	7.66 ± 0.17	3.53 ± 0.07
MixMatch (Berthelot et al., 2019)	WRN-28-2	1.5M	6.24 ± 0.06	2.89 ± 0.06

Methods	SSL	10%	100%
ResNet-50	✗	55.09 / 77.26	77.28 / 93.73
w. RandAugment		58.84 / 80.56	78.43 / 94.37
UDA (RandAugment)	✓	68.78 / 88.80	79.05 / 94.49

Table 5: Top-1 / top-5 accuracy on ImageNet with 10% and 100% of the labeled set. We use image size 224 and 331 for the 10% and 100% experiments respectively.

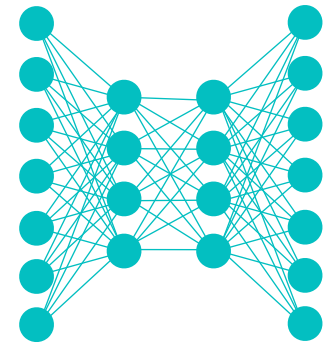




Lecture Notes for **Neural Networks and Machine Learning**

SSL

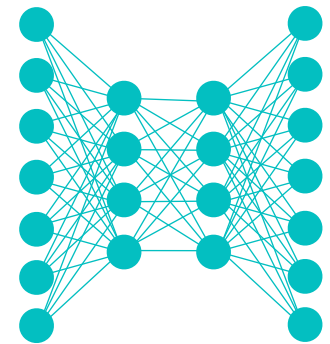
Next Time:
MML and MTL
Reading: None



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



Multi-task and
Multi-Modal Learning



Logistics and Agenda

- Logistics
 - Grading Update
- Agenda
 - Student Paper Presentation
 - Multi-modal and Multi-task
- Next Time
 - Multi-task demo and Town Hall
 - Finish Demos



Multi-modal Review



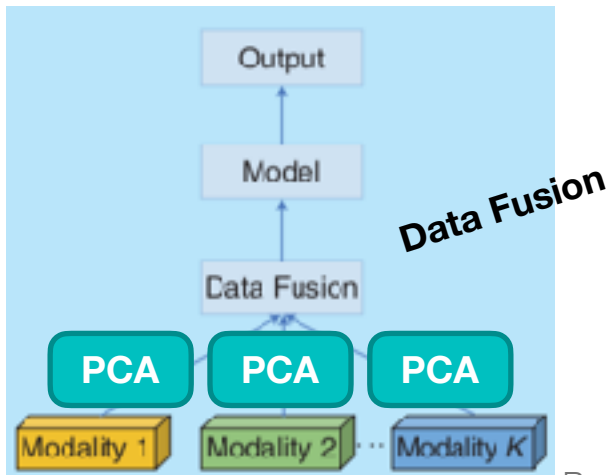
Multi-modal == Multiple Data Sources

- **Modal** comes from the “sensor fusion” definition from Lahat, Adali, and Jutten (2015) for deep learning
- Using the Keras functional API, this is extremely easy to implement
 - ... and we have used it since CS7324!
- But now let's take a deeper dive and ask:
 - What are the different types of modalities that we might try?
 - Is there a more optimal way to merge information?
 - When? Early, Intermediate, and late fusion



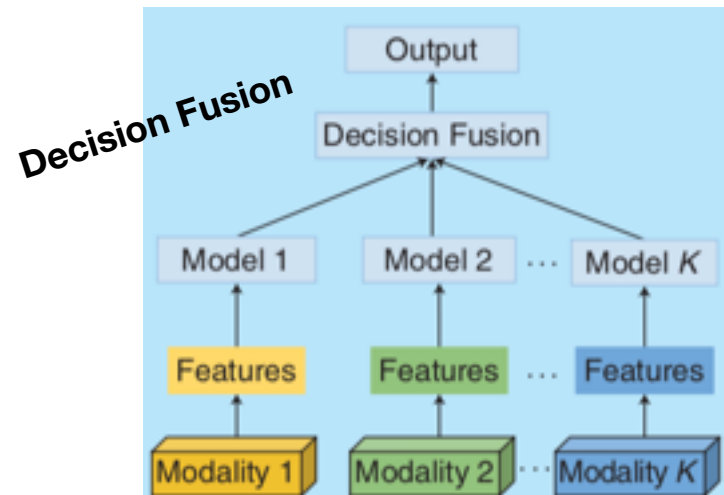
Early and Late Stage Fusion

- **Early Fusion:** Merge sensor layers early in the process
- **Assumption:** there is some data redundancy, but modes are conditionally dependent
- **Problem:** architecture parameter explosion
 - Typically need dimensionality reduction



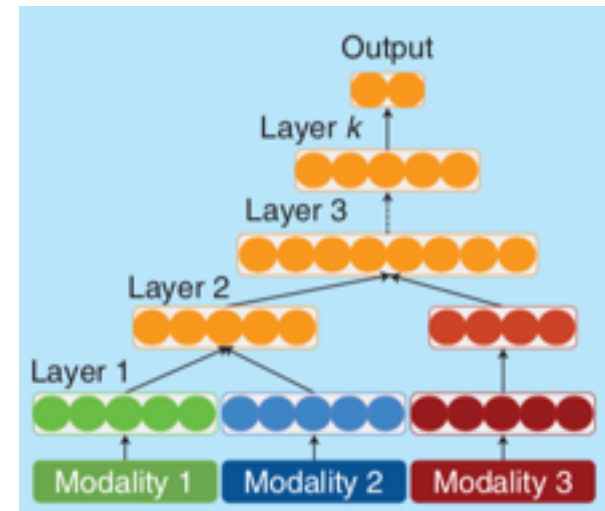
Ramamchandran and Taylor, 2017

- **Late Fusion:** Merge sensor layers right before flattening
- Use Decision Fusion on outputs
- **Assumption:** little redundancy or conditional independence—just an ensemble architecture
- **Problem:** just separate classifiers, limited interplay



Intermediate Fusion

- Merge sensor layers in soft way
 - **Assumption:** some features interplay and others do not
 - **Problem:** how to optimally tie layers together?
1. Stacked Auto-Encoders [Ding and Tao, 2015]
 2. Early fuse layers that are correlated [Neverova *et al.* 2016]
 3. Fully train each modality merge based on criterion of similarity in activations [Lu and Xu 2018]
 4. Granger Cluster data in each modality and combine [Sylvester *et al.* 2023]



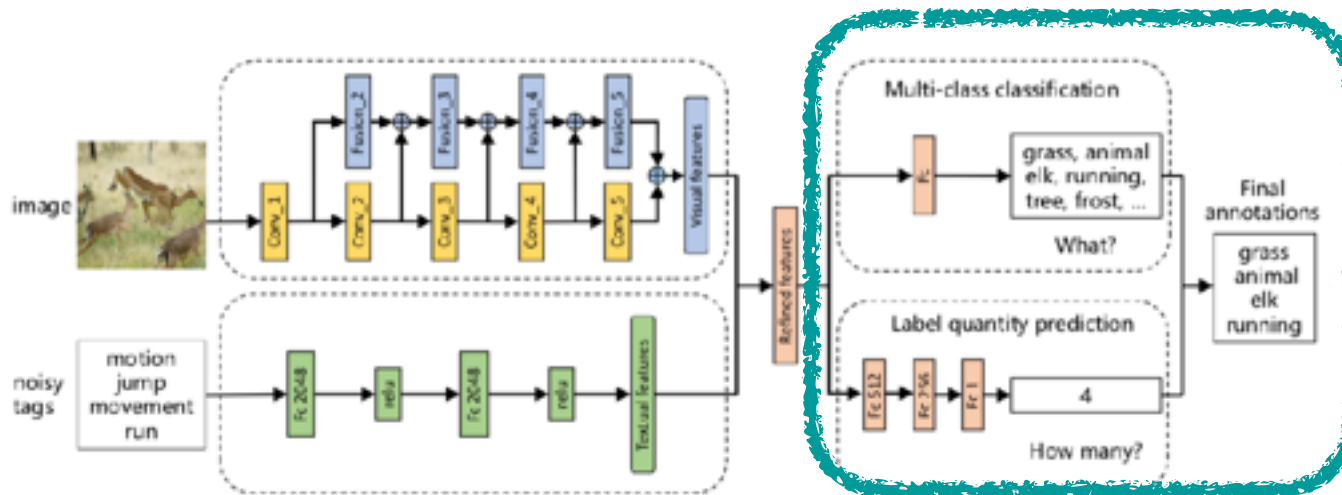
Ramamchandran and Taylor, 2017

103



Multi-modal Merging

- **Still an open research problem**
- How to develop merging techniques that
 - Can handle exponentially many pairs of modalities
 - Automatically merge meaningful modes
 - Discard poor pairings
 - Selectively merge early or late (or dynamically)



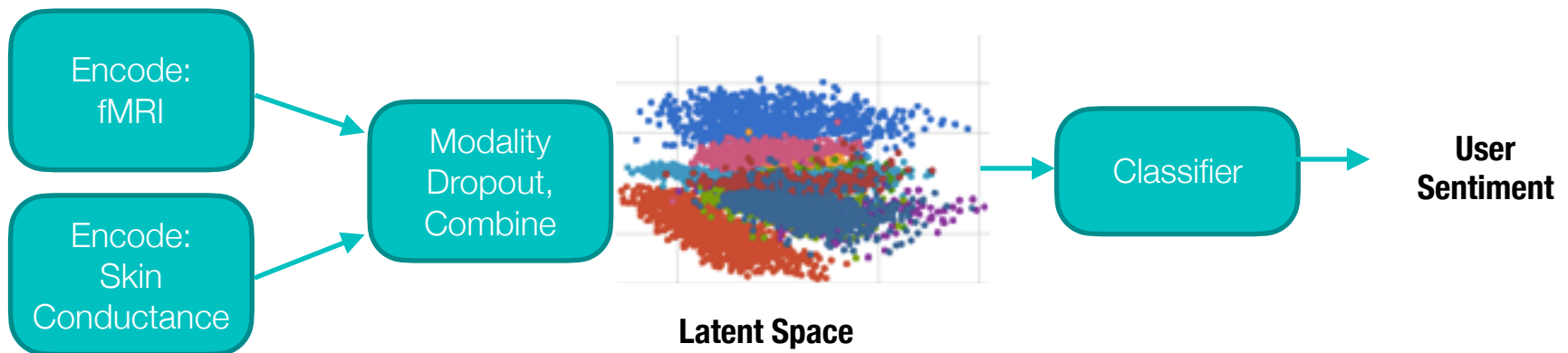
<https://arxiv.org/pdf/1709.01220.pdf>

Most current methods are still ad-hoc



Approaches with Deep Learning

- Latent Space Transfer (universality)
 - From another domain, map to a similar latent space for the same task
 - Useful for unifying data based upon a new input mode when old mode is well understood
 - ◆ for example, biometric data
 - ◆ **2019-2023, I have never seen a research paper on this...**



High-Modality Multimodal Transformer: Quantifying Modality & Interaction Heterogeneity for High-Modality Representation Learning

Paul Pu Liang¹, Yiwei Lyu², Xiang Fan¹, Jeffrey Tsaw¹, Yudong Liu¹, Shentong Mo¹, Dani Yogatama³, Louis-Philippe Morency¹, Ruslan Salakhutdinov¹

¹Carnegie Mellon University, ²University of Michigan, ³DeepMind

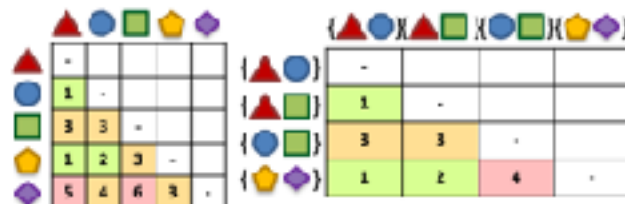
1a. Estimate modality heterogeneity via transfer



1b. Estimate interaction heterogeneity via transfer



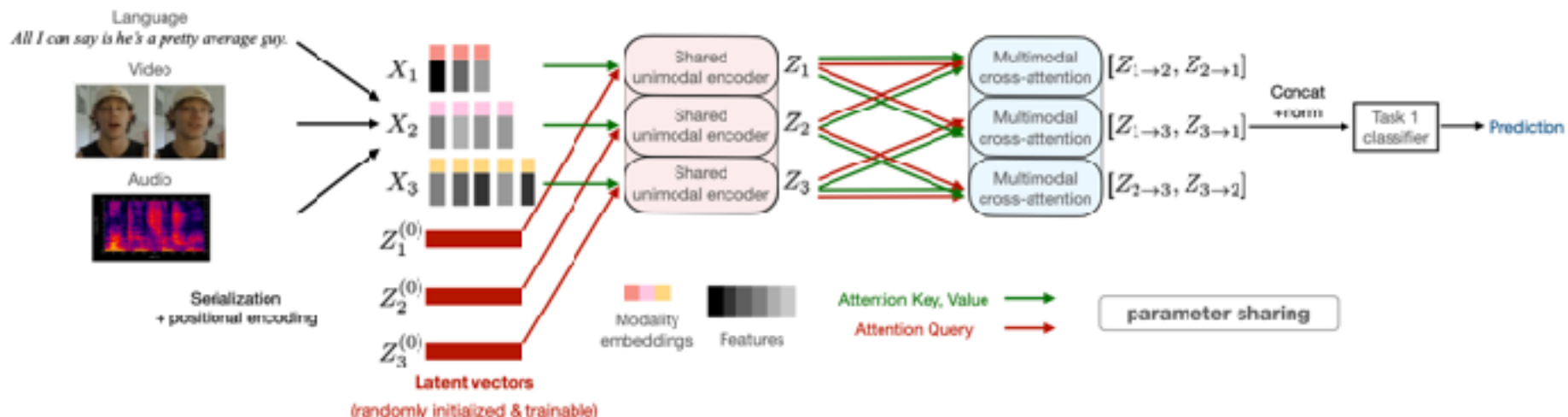
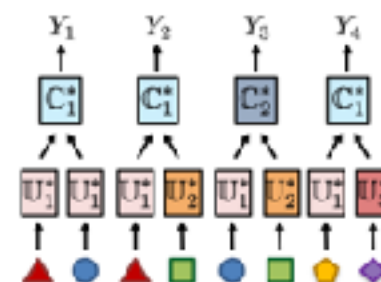
2a. Compute modality & interaction heterogeneity matrices



2b. Determine parameter clustering

$$\begin{aligned} U_1 &= \{U_1, U_2, U_4\} & C_1 &= \{C_{12}, C_{13}, C_{45}\} \\ U_2 &= \{U_3\} & C_2 &= \{C_{23}\} \\ U_3 &= \{U_5\} \end{aligned}$$

3. Heterogeneity-aware model across modalities and tasks

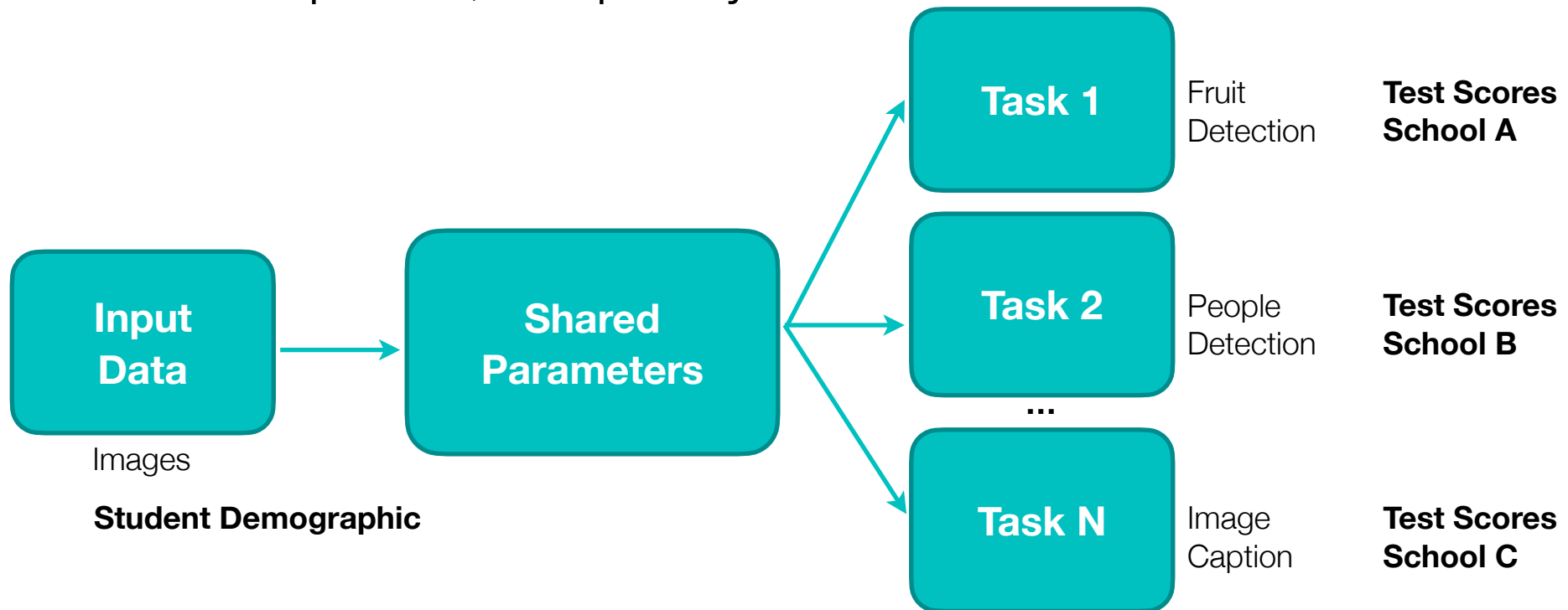


Multi-Task Models

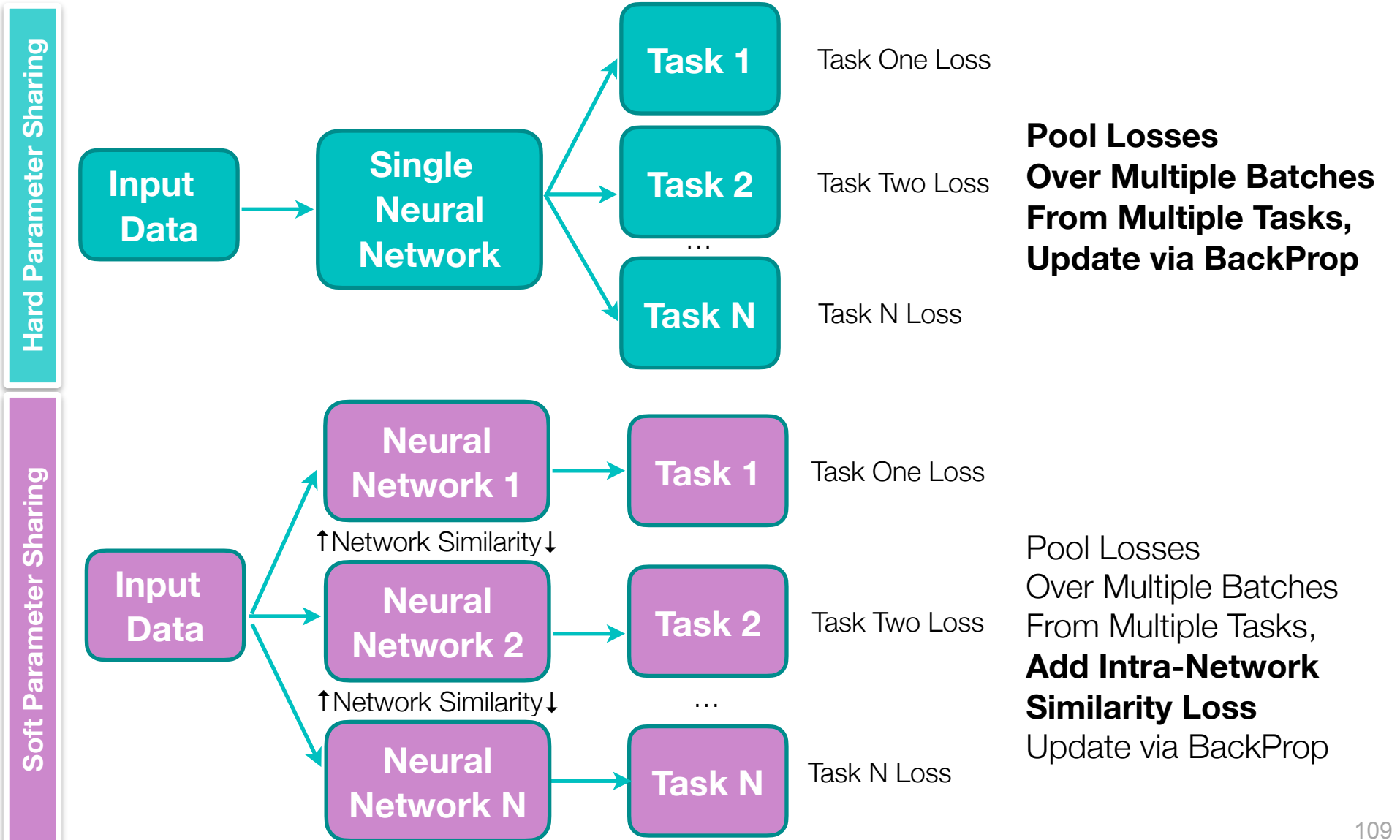


Multi-task learning overview

- For deep networks, simple idea: share parameters in early layers
- Used shared parameters as feature extractors
- Train separate, unique layers for each task

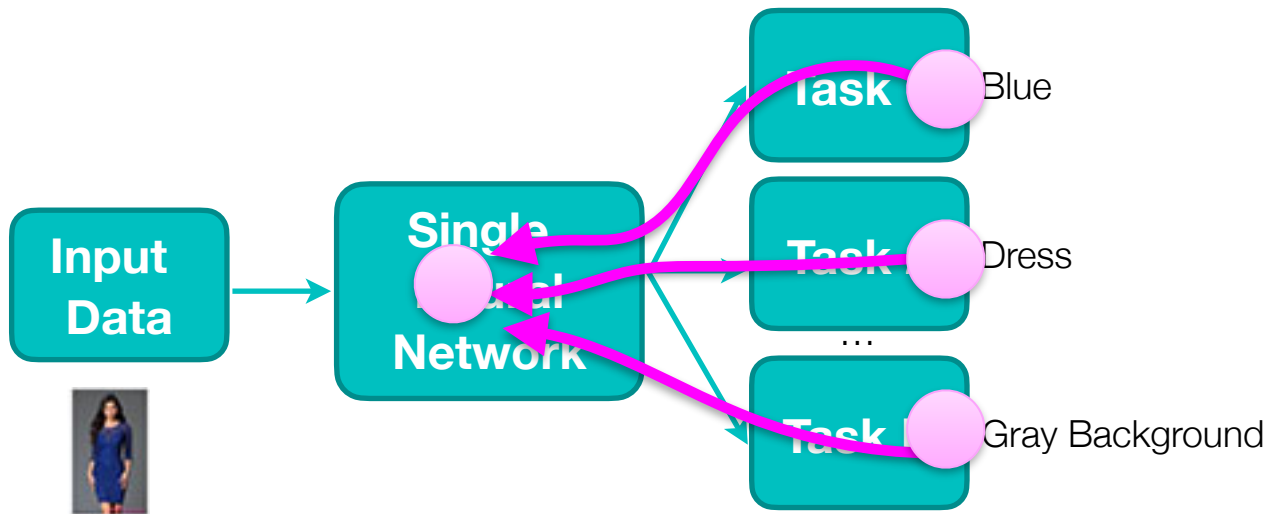


Multi-task Learning Parameter Sharing



Multi-task Optimization

Multi-Label per Input



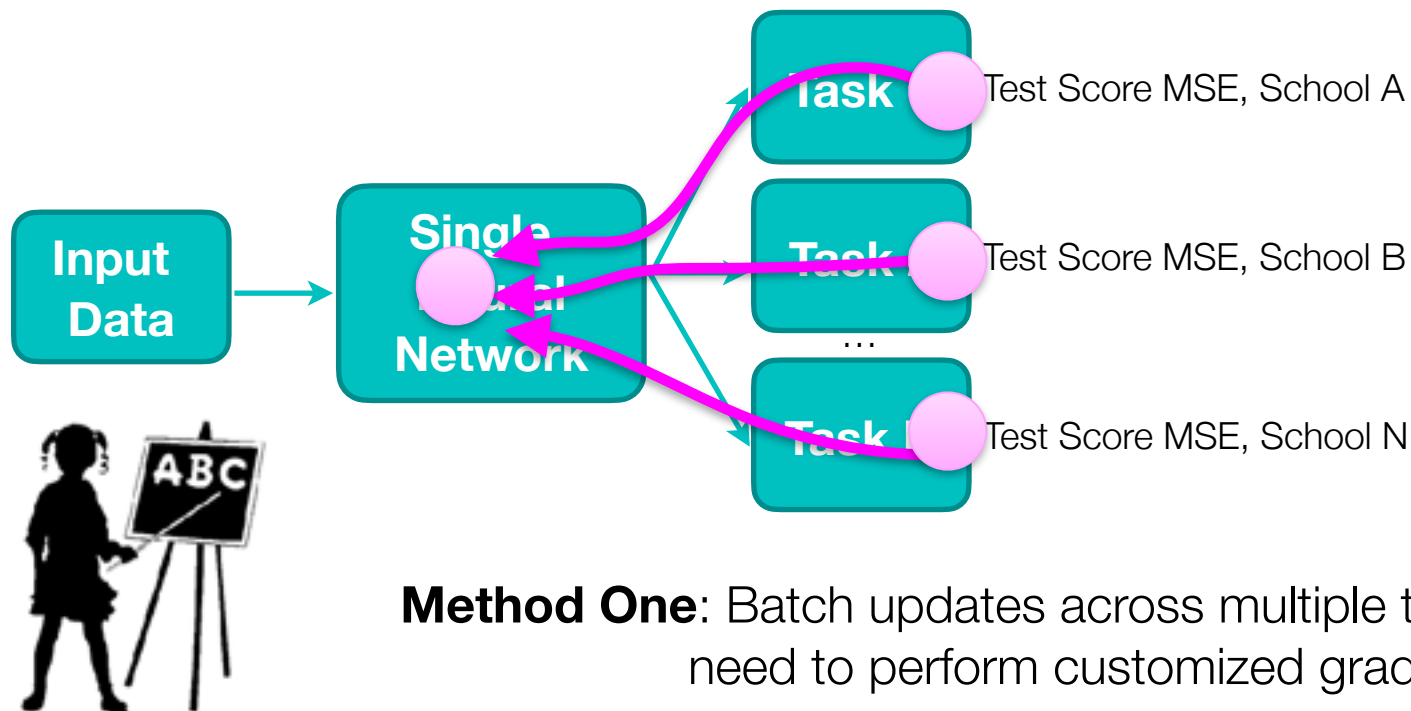
Measure Loss **for each label simultaneously**

Back propagate **everything at one time** for a given batch



Multi-task Optimization

Single Task Label per Input



- Method One:** Batch updates across multiple tasks
need to perform customized gradient calculations
- Method Two:** Update small batches using a random task
easier, but can cause instability in training





Optional Demo

Multi-Task Learning in Keras with **Multi-Label Data**

Fashion week, colors and dresses

Follow Along: <https://www.pyimagesearch.com/2018/06/04/keras-multiple-outputs-and-multiple-losses/>





Multi-Task Learning

School Data, Computer Surveys



Traian-Pop Traian Pop



LukeWood Luke Wood

KerasCV Author, Full Time Keras team member & Machine Learning researcher @ Google, Part Time UCSD Ph.D student



♡ Sponsor

Follow

Method One: Batch updates across multiple tasks
need to perform customized gradient calculations

Method Two: Update small batches using a random task
easier, but can cause instability in training

Follow Along: [LectureNotesMaster/03](#) [LectureMultiTask.ipynb](#)

113



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

Multi-Modal and Multi-Task



Next Time:
Circuits

Reading: Chollet 8.1-8.5

