

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



Multi-task Learning



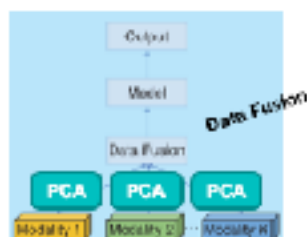
Logistics and Agenda

- Logistics
 - Back to regular office hours
 - ◆ No virtual option?
- Agenda
 - Multi-task networks and demo
 - Town Hall
- Next Time
 - CNN Visualization



Last Time

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



Decision Fusion and Textor, 2017



90

Neural Architecture Search for Mode Fusion

Genetic Algorithm

1. Sample new candidates
2. Evaluate fitness
3. Mutate and Crossover
4. Keep the best solutions
5. Repeat

Very computational when starting, because candidates are all untrained. However, as more blocks start from "mostly trained" positions, training time reduces.

1. Sample from set of Random
2. Train by fine-tuning
3. Evaluate fitness
4. Add fitness function back to "Possible Fusions"

Found solution for AW-MNIST

Applied Intelligence
<https://doi.org/10.1007/s10489-022-04085-x>

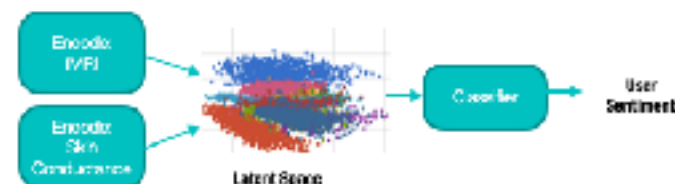


An approach for combining multimodal fusion and neural architecture search applied to knowledge tracing

Xinyi Ding¹ · Tao Han¹ · Yili Fang¹ · Eric Larson²

Accepted: 18 August 2022
 © The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2022

- 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
- I have never seen a research paper on this...



98

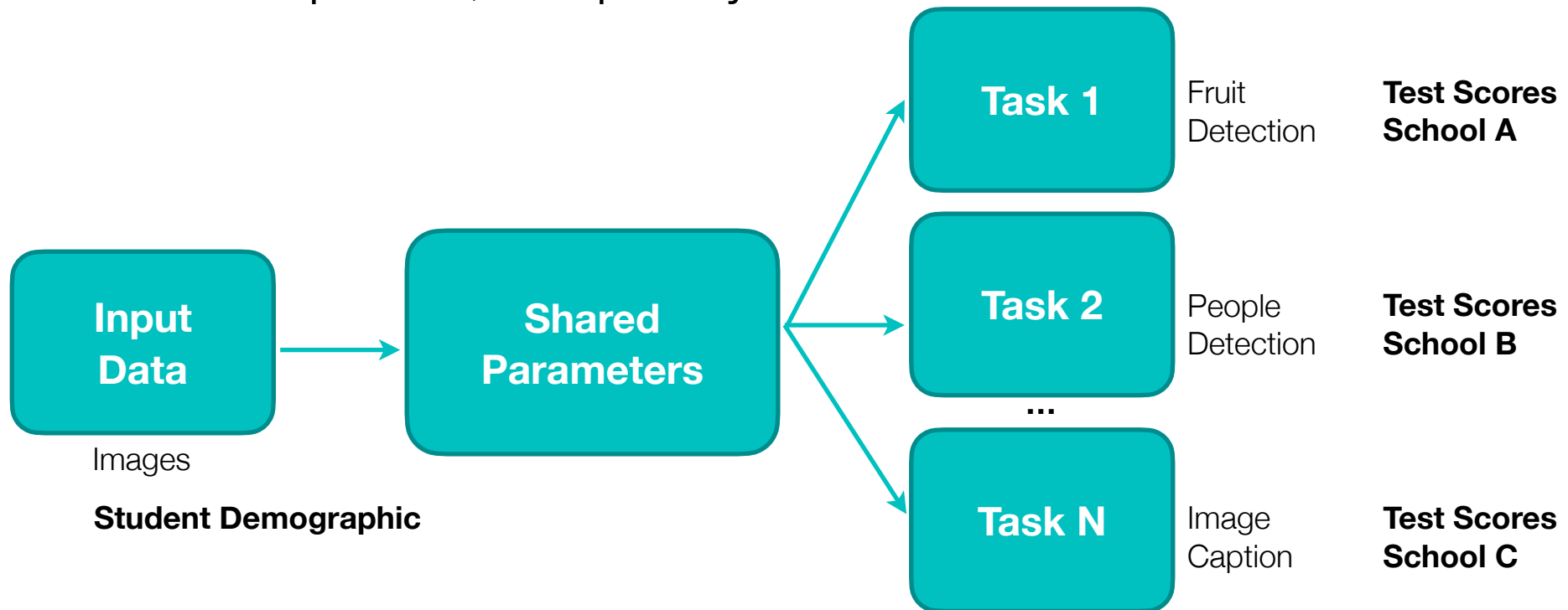


Multi-Task Models

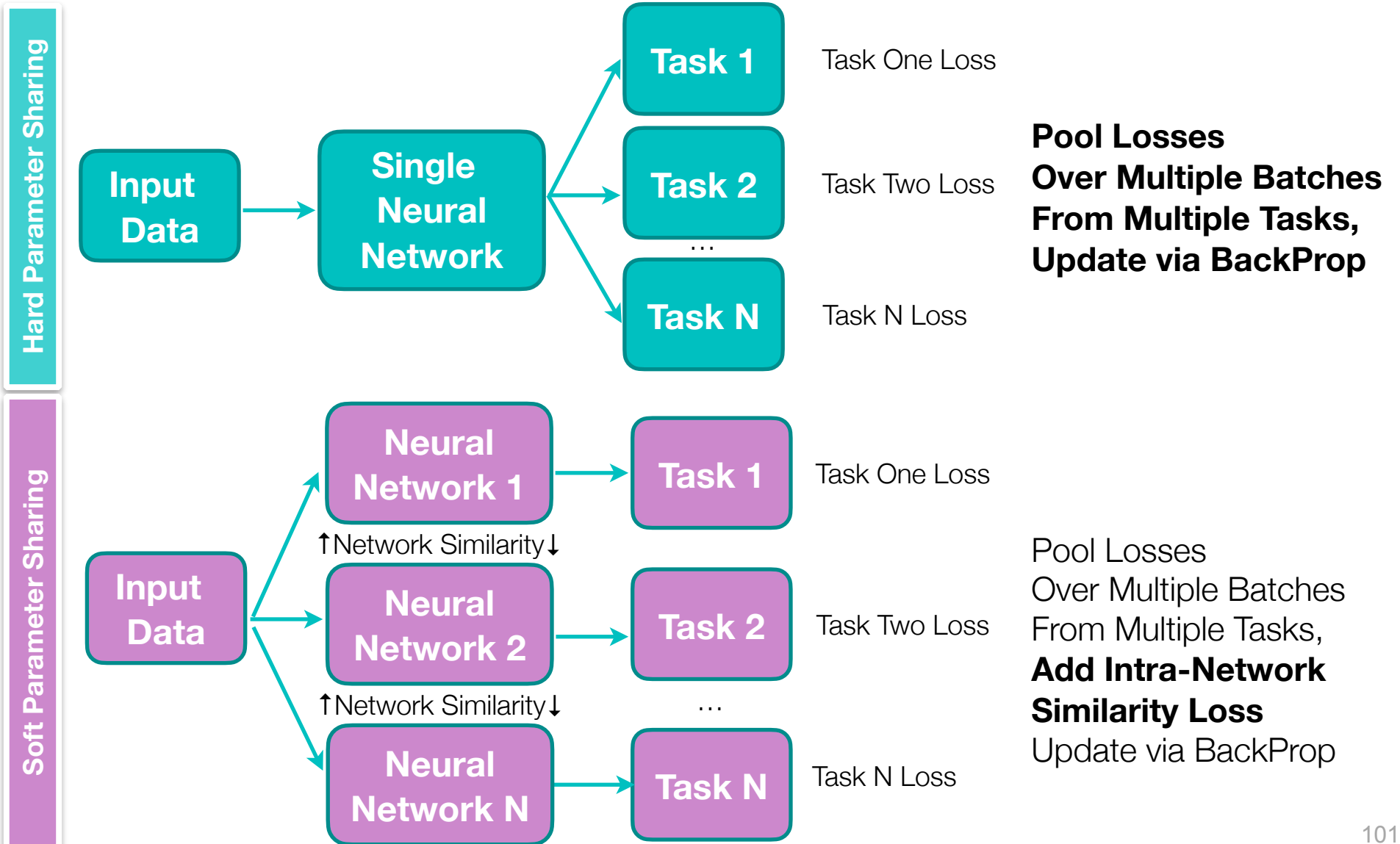


Multi-task learning overview

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

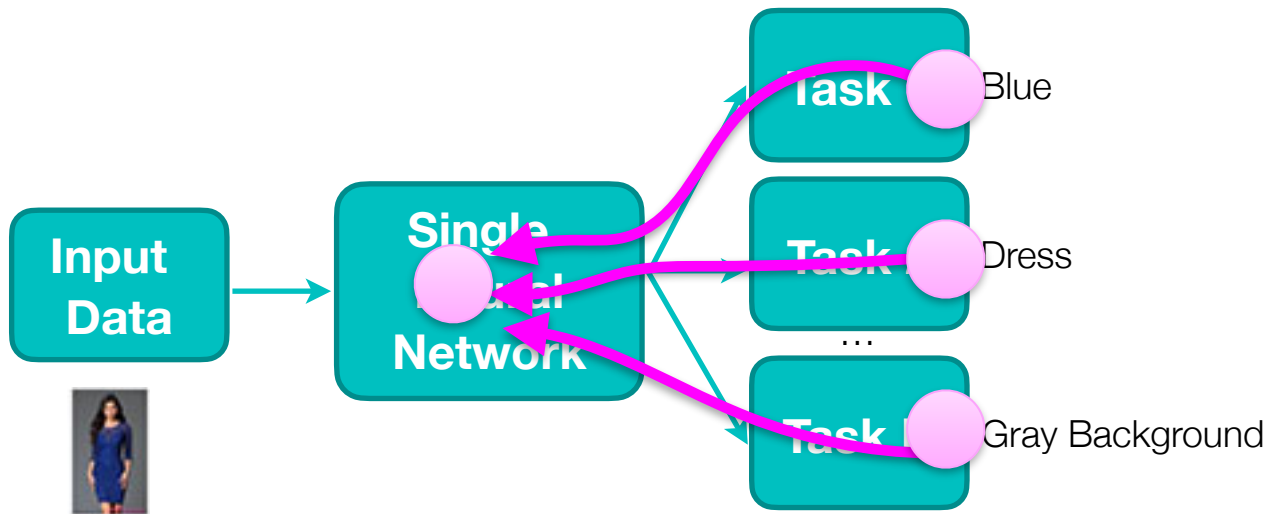


Multi-task Learning Parameter Sharing



Multi-task Optimization

Multi-Label per Input



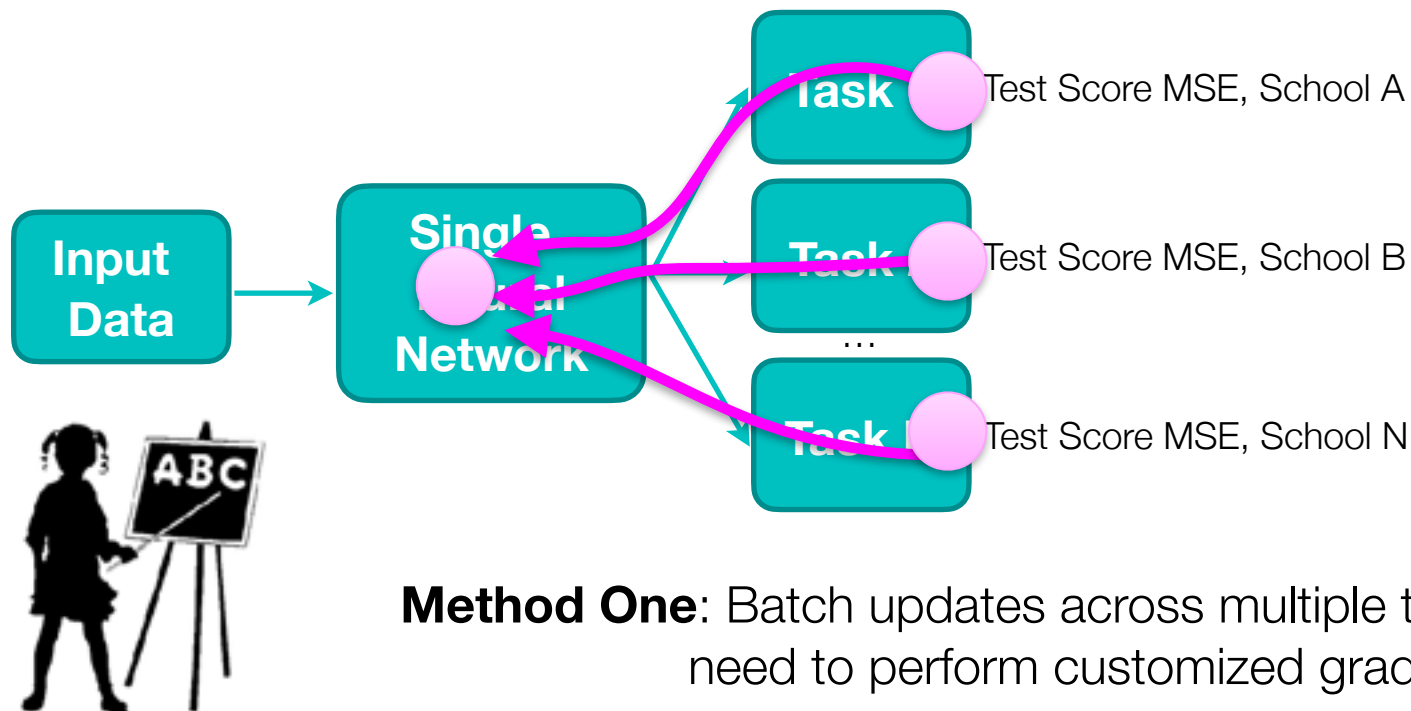
Measure Loss **for each label simultaneously**

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



Follow

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

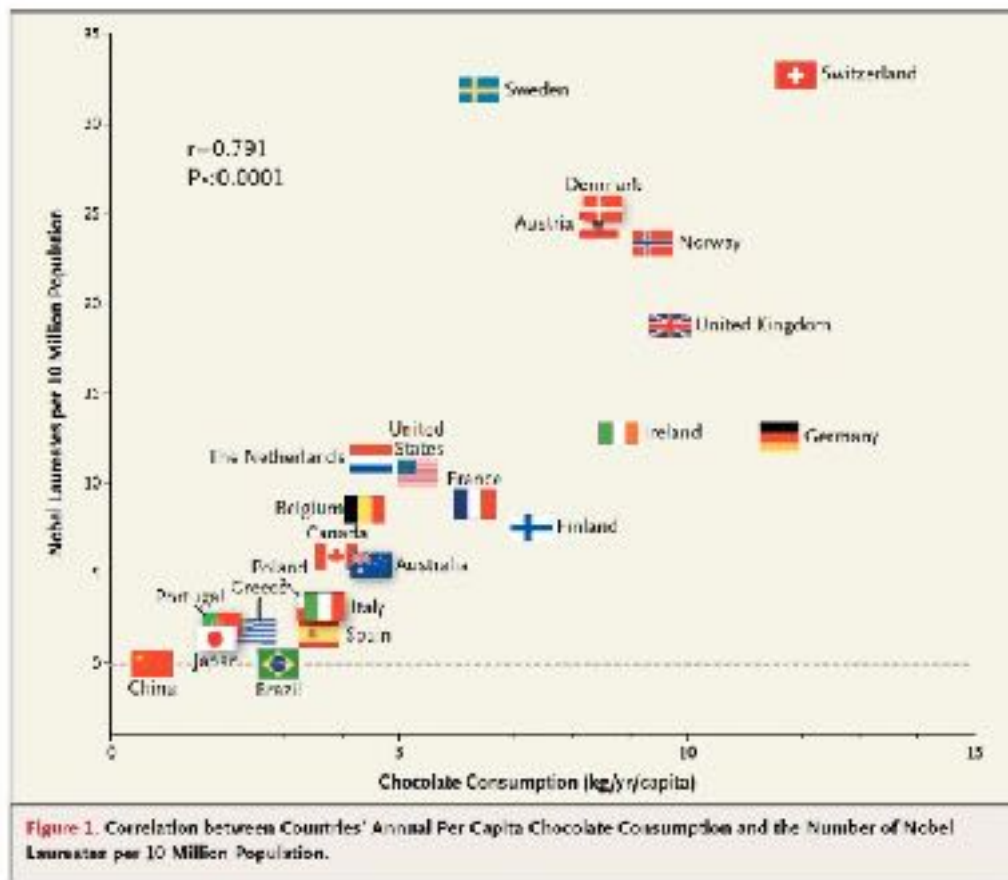
Method Two: Update small batches using a random task
easier, but can cause instability in training

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

105



Lab Two Town Hall



Multi-Task Networks
Multi-Modal Networks



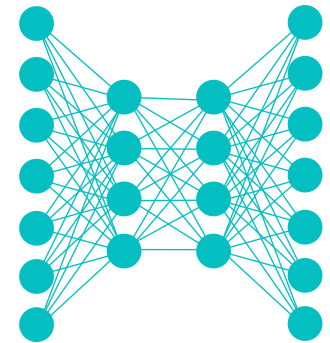
Lecture Notes for **Neural Networks and Machine Learning**

Multi-Modal and Multi-Task



Next Time:
Circuits

Reading: Chollet 8.1-8.5





Backup Slides



Active Transfer Learning

Theory:



Practice:



Machine Learning :

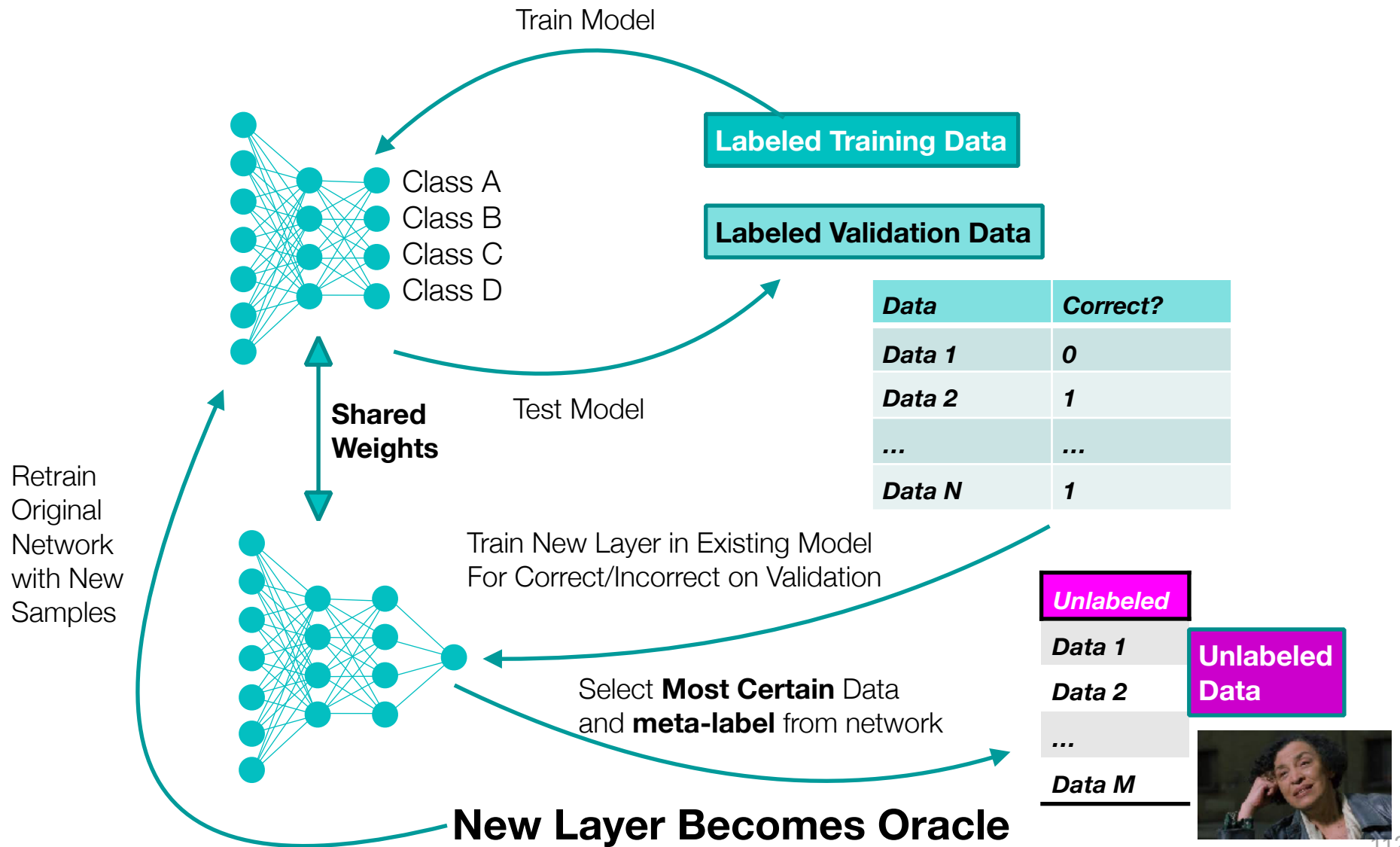


Active Learning Overview

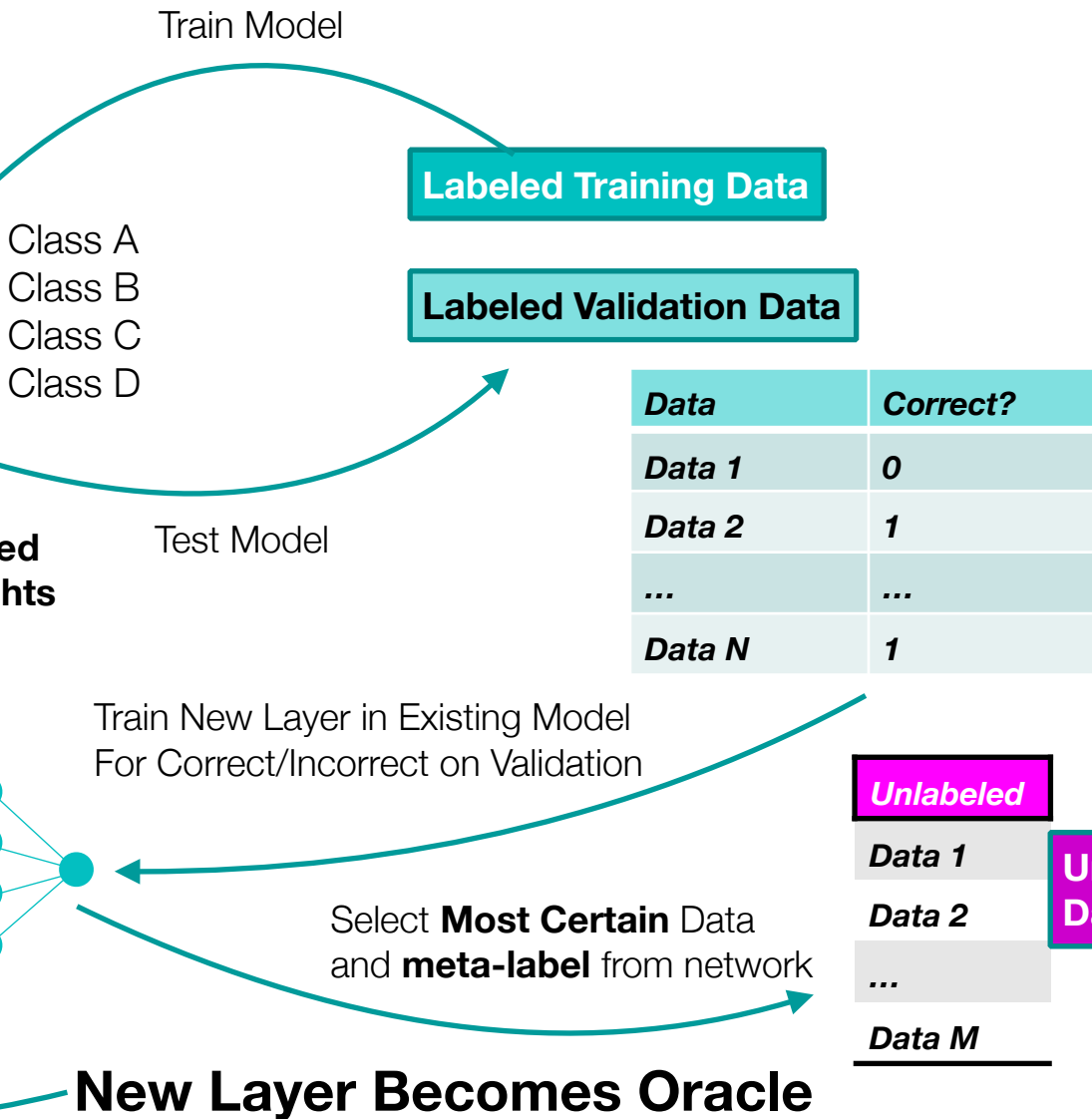
- **Basic Idea:** Use a trained model to sample from an oracle that can magically give you a new label
 - Active Learning:
What labels should we ask the oracle about?
- Uncertainty Sampling
 - Choose instances where the model is most uncertain or most certain
 - Various ways to measure certainty
- Diversity Sampling
 - Choose instances that are similar or different from training distribution



Uncertainty Sampling with a Neural Network



Uncertainty Sampling with a Neural Network

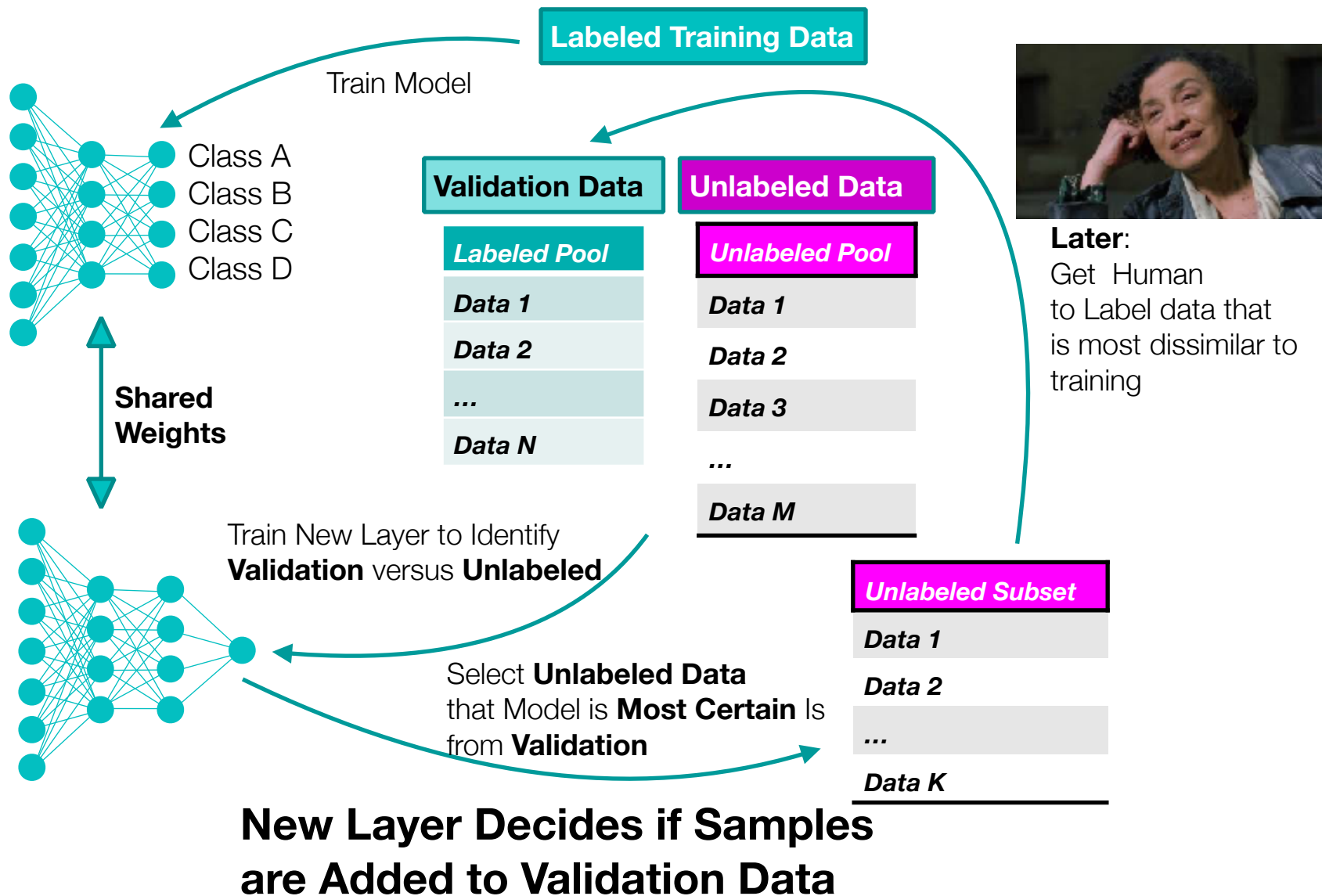


Problems:

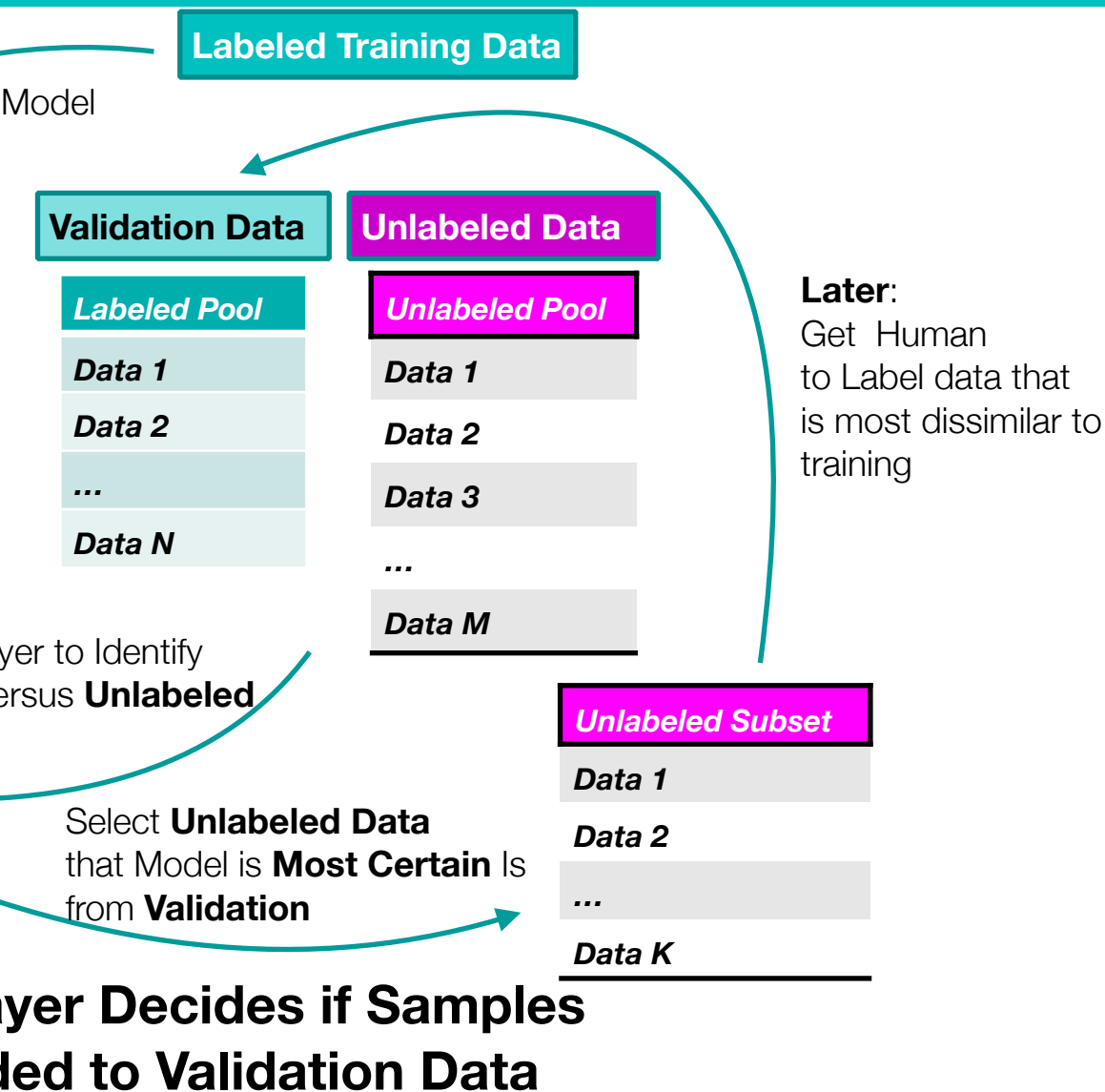
- Training pool is represented by classes the model already does well predicting
- Limited diversity of Samples
- Training pool can become contaminated easily from a few wrong predictions
- For Oracle: we might be asking to get labels that the model is already good at classifying



Diversity Sampling with a Neural Network



Diversity Sampling with a Neural Network

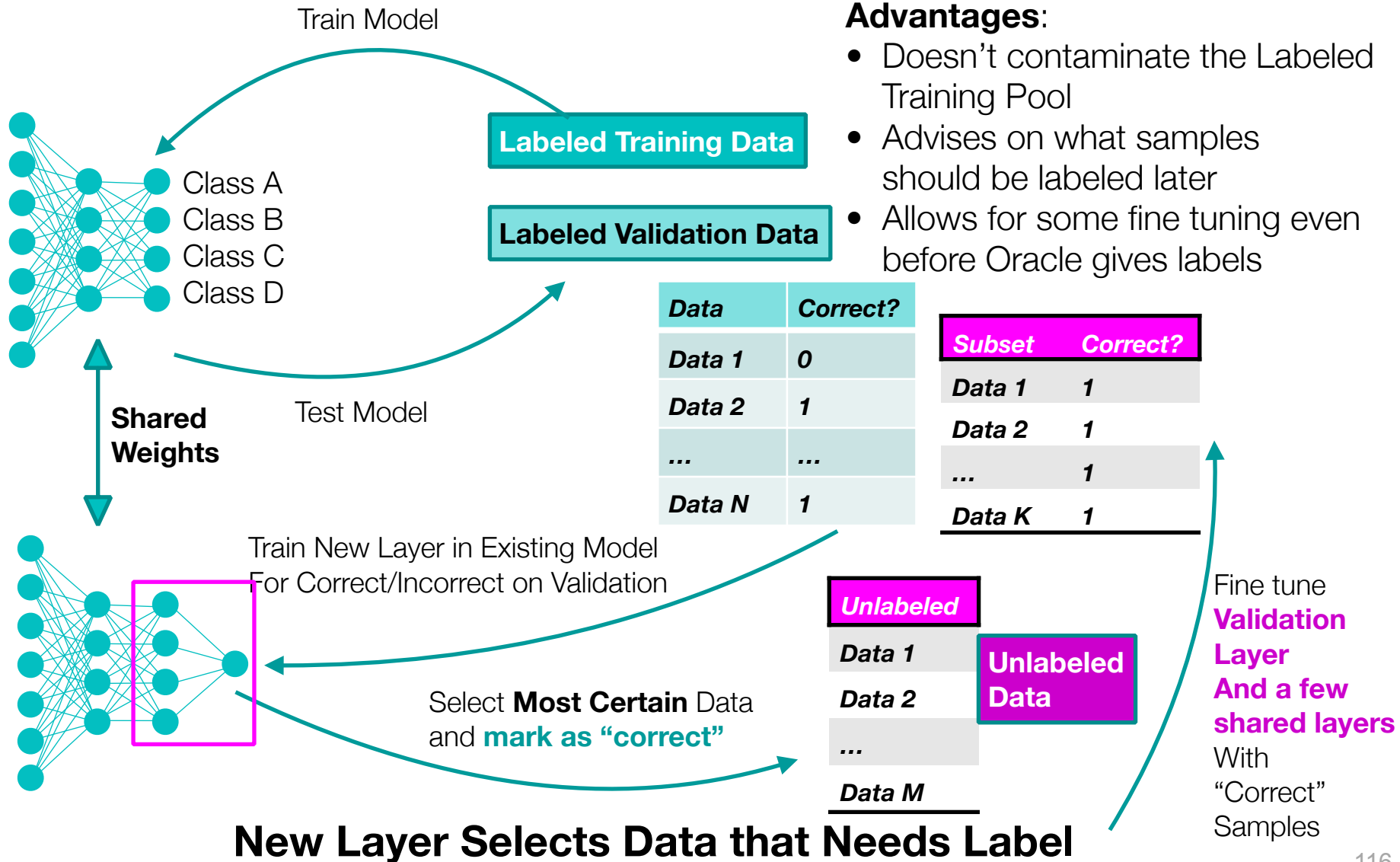


Discussion:

- Training pool is not contaminated
- Expands validation data in well mannered way, not adding too “far away” samples
- Validation versus Unlabeled might not be the best comparison, because it ignores confusions in the training data
- For Oracle: we can get labels to inputs that the model is likely to be unsure about
- But... this only helps us when we have an Oracle to give us labels



ATLAS: Active Transfer Learning for Adaptive Sampling



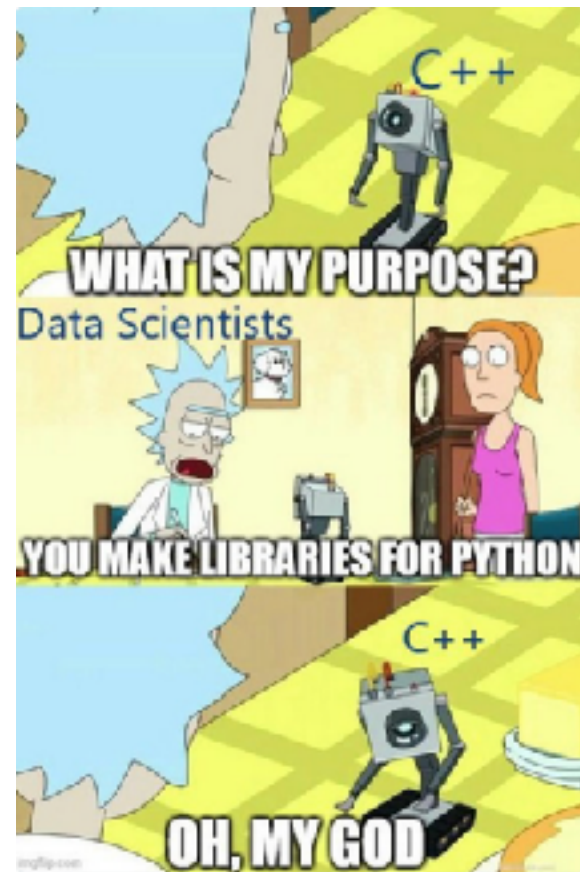
Time Period	Protocol	Expected Feedback
First Week	<p>Homeowner provides 8-20 examples over the first week:</p> <ul style="list-style-type: none"> • 1-2 Shower usages • 1 run of the dishwasher • 1 run of the laundry machine • 2 examples of each toilet • 1 example of hot and cold water use for each dual handle faucet • 1 example of hot, cold, and mixed water use for each single handle faucet (2 examples if in kitchen) 	<p>HydroSense relies on the rule based classifier for the first week.</p> <p>Pressure waves are saved in order to create a sparse codebook of features.</p> <p>Results are displayed at the fixture category for dishwashers, showers, and washing machines.</p>
Start of Second Week	Homeowner provides 2-4 labels every other day when the system messages them on their mobile device	<p>Results are displayed at the full fixture category level from the CoDBN-VE algorithm. Expected accuracy:</p> <ul style="list-style-type: none"> • 85% at fixture category level
End of Second Week	Homeowner has supplied 9-12 examples that were flagged by active learning.	<p>HydroSense now displays results at the Lumped Fixture level.</p> <p>Expected accuracy:</p> <ul style="list-style-type: none"> • 82% at fixture level • 87% at fixture category level
End of Third Week	Homeowner continues to supply sparsely selected examples every other day. About 9-12 additional examples provided.	<p>Valve level accuracy now provided.</p> <p>Expected accuracy:</p> <ul style="list-style-type: none"> • 80% at valve level • 87% at fixture level • 92% at fixture category level
Fourth Week	Homeowner can optionally continue to provide examples to the system for increased accuracy.	<p>Expected accuracy:</p> <ul style="list-style-type: none"> • 81% at valve level • 89% at fixture level • 93% at fixture category level

Table 8-2. Expected feedback and calibration protocol for semi-supervised HydroSense system

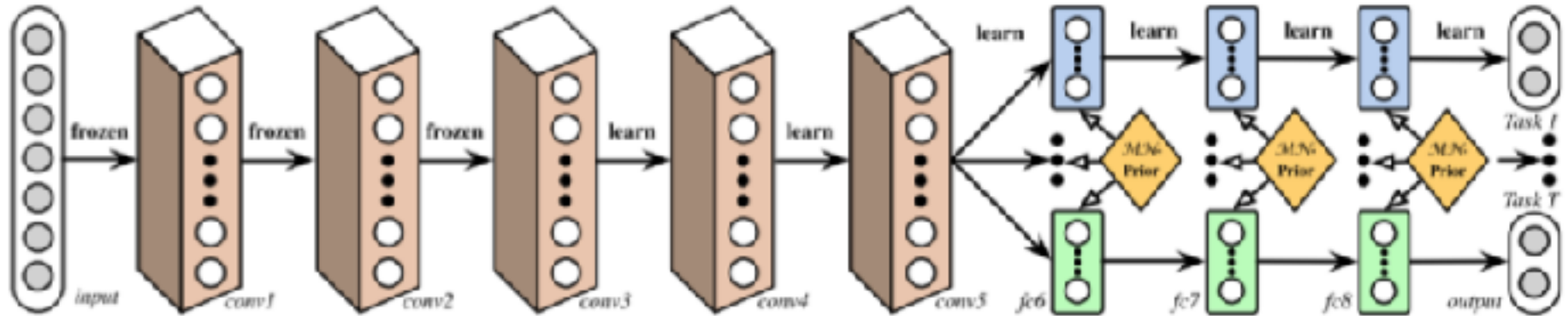
Optional Topic: Multi-Task Model Examples

He uses statistics like a drunken man uses a lamp post, more for support than illumination.

-- Andrew Lang



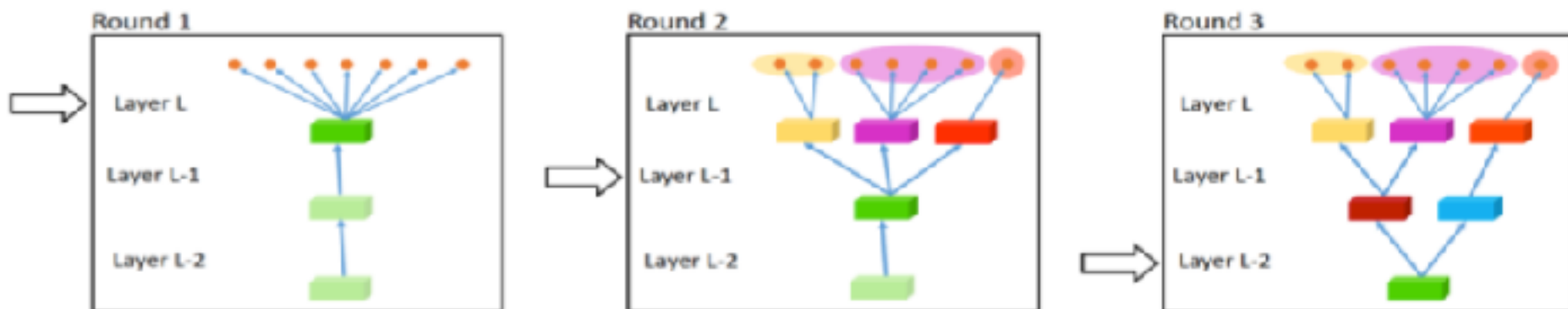
Multi-task: Deep Relationship Networks



- Start training traditionally (CCE)
- Minimize Kroenecker Product of fully connected task specific layers (here matrices are vectorized and therefore it is an outer product)
 - intuitively: make Covariances between tasks close to a given prototype Covariance
 - encourages feature maps in each task to be **less correlated** to feature maps of another task



Multi-task: Adaptive Feature Sharing



- Train
- Rep

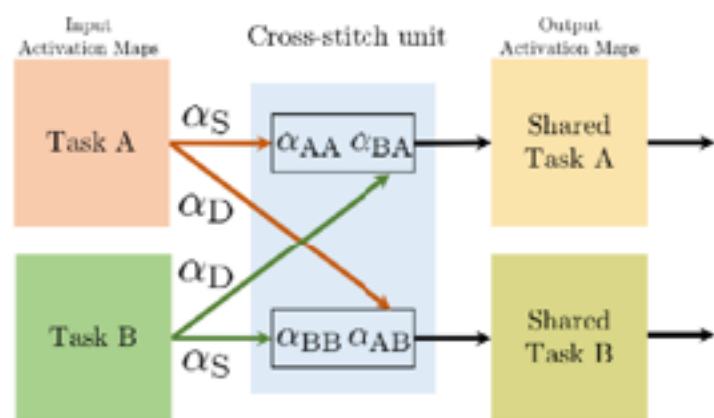
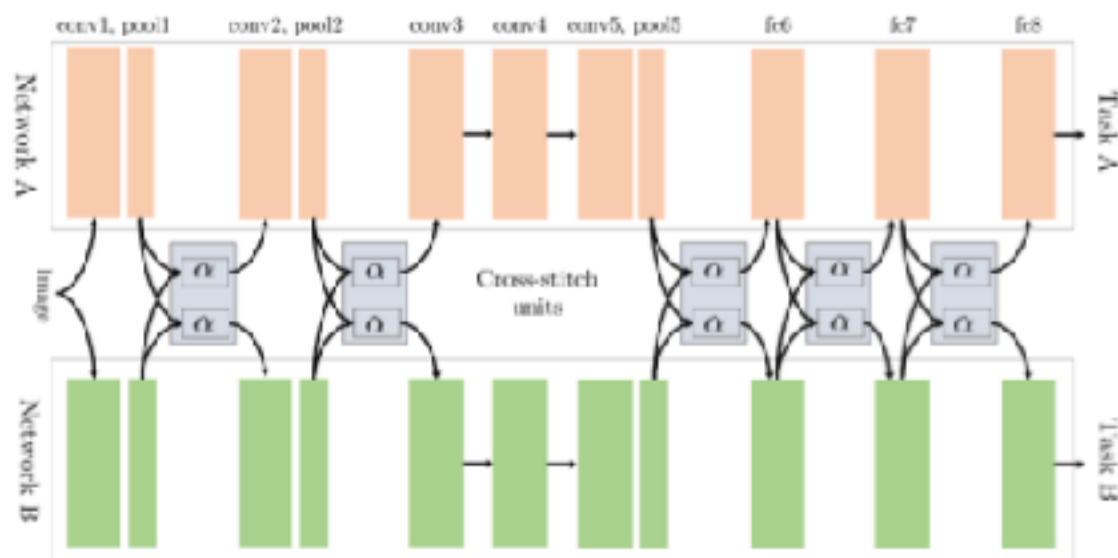
$$A^*, \omega^*(l) = \arg \min_{A \in \mathbb{R}^{d \times d'}, |\omega| = d'} ||W^{p,l} - AW_{\omega}^{p,l}||_F, \quad (2)$$

where $W_{\omega}^{p,l}$ is a truncated weight matrix that only keeps the rows indexed by the set ω . This problem is NP-hard, however, there exist approaches based on convex relaxation

- Cluster affinity of branch is not final layer
- Cut weights and retrain (fine tune) network
- Decrement current layer index



Multi-task: Cross Stitch Networks

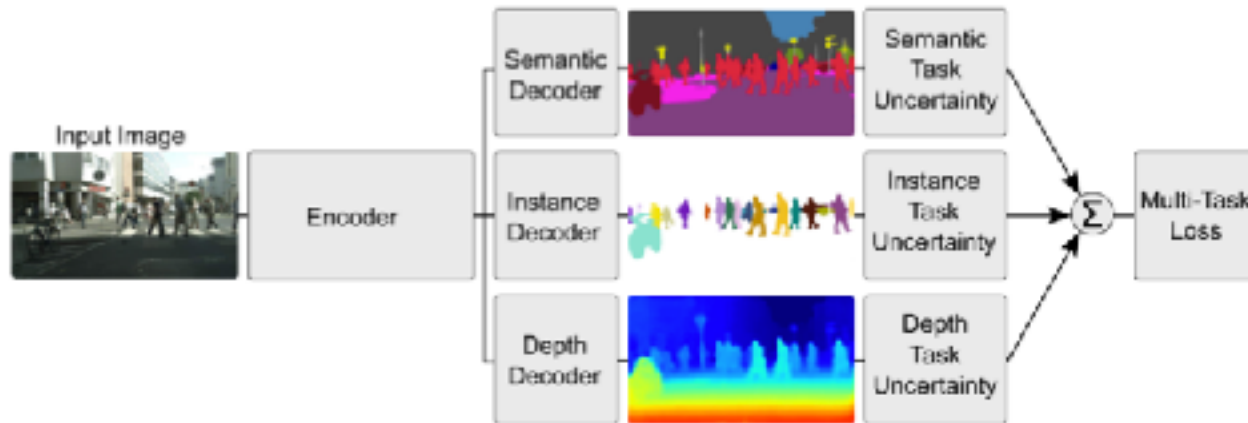


$$\begin{bmatrix} \hat{x}_A^{ij} \\ \hat{x}_B^{ij} \end{bmatrix} = \begin{bmatrix} \alpha_{AA} & \alpha_{AB} \\ \alpha_{BA} & \alpha_{BB} \end{bmatrix} \begin{bmatrix} x_A^{ij} \\ x_B^{ij} \end{bmatrix}$$

- Only works for simultaneous multi-label problems
 - like semantic segmentation and surface normal segmentation (clustering similarly facing objects)
- Take a learned weighted sum of the activations
- Works a little better than single task, but no worse



Multi-task: Uncertainty Weighting



- Use variance of each loss function from each task to normalize
 - Which is a great idea, tasks with more uncertainty become less influential to loss function
 - call it homoscedastic without sound reasoning because that feels better than “normalized variance”
 - ◆ talk about homoscedasticity for no reason
- Write an entire paper in a “mathy” way to make contribution less clear
- Reviewers give you a pass because you have a history of doing good work



Current Multi-task Research

- Incredibly diverse sets of solutions
- Mostly not evaluated on similar datasets
- Reasoning given is mostly ad-hoc...
- Theory is wildly under developed
 - because the problem is incredibly difficult
- Neural architecture search is an option...



AlphaCode

