

Machine Learning Algorithms for Low Data Environments

Alan Lu, Department of Electrical and Computer Engineering at the University of Illinois
at Urbana-Champaign

Mentors: Neil Getty

Supervisor: Fangfang Xia

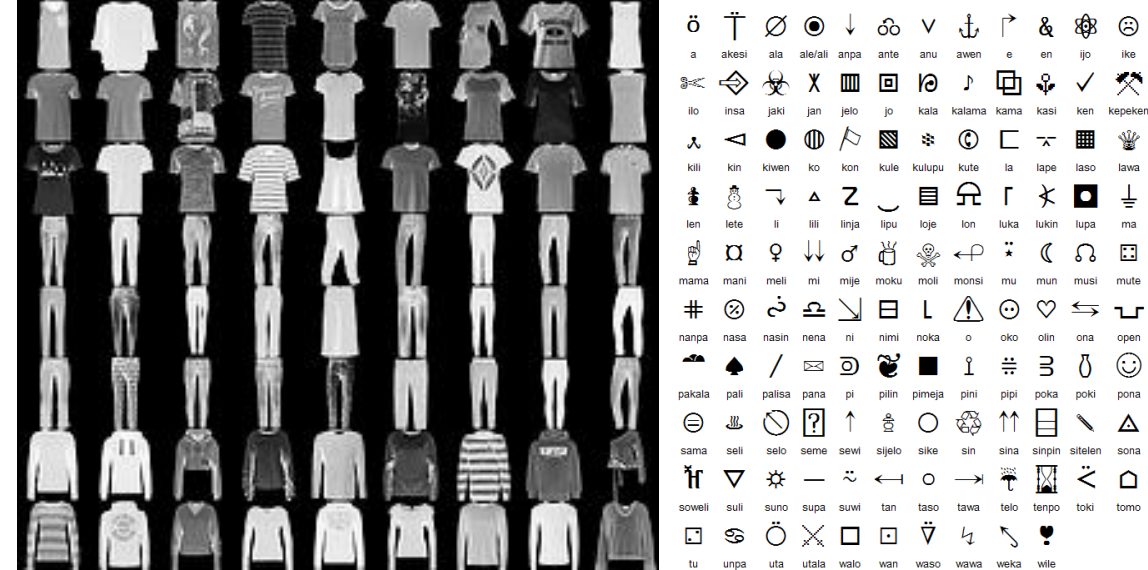


Motivation

- ▶ Given a large amount of data, current state-of-the-art machine learning algorithms perform with high accuracy
- ▶ In many industries obtaining data is expensive and difficult
- ▶ It is necessary to apply machine learning to these industries with small data
 - ▶ Cancer patient data
 - ▶ Pharmaceutical applications
 - ▶ One-shot learning applications - image recognition, drug discovery

Datasets

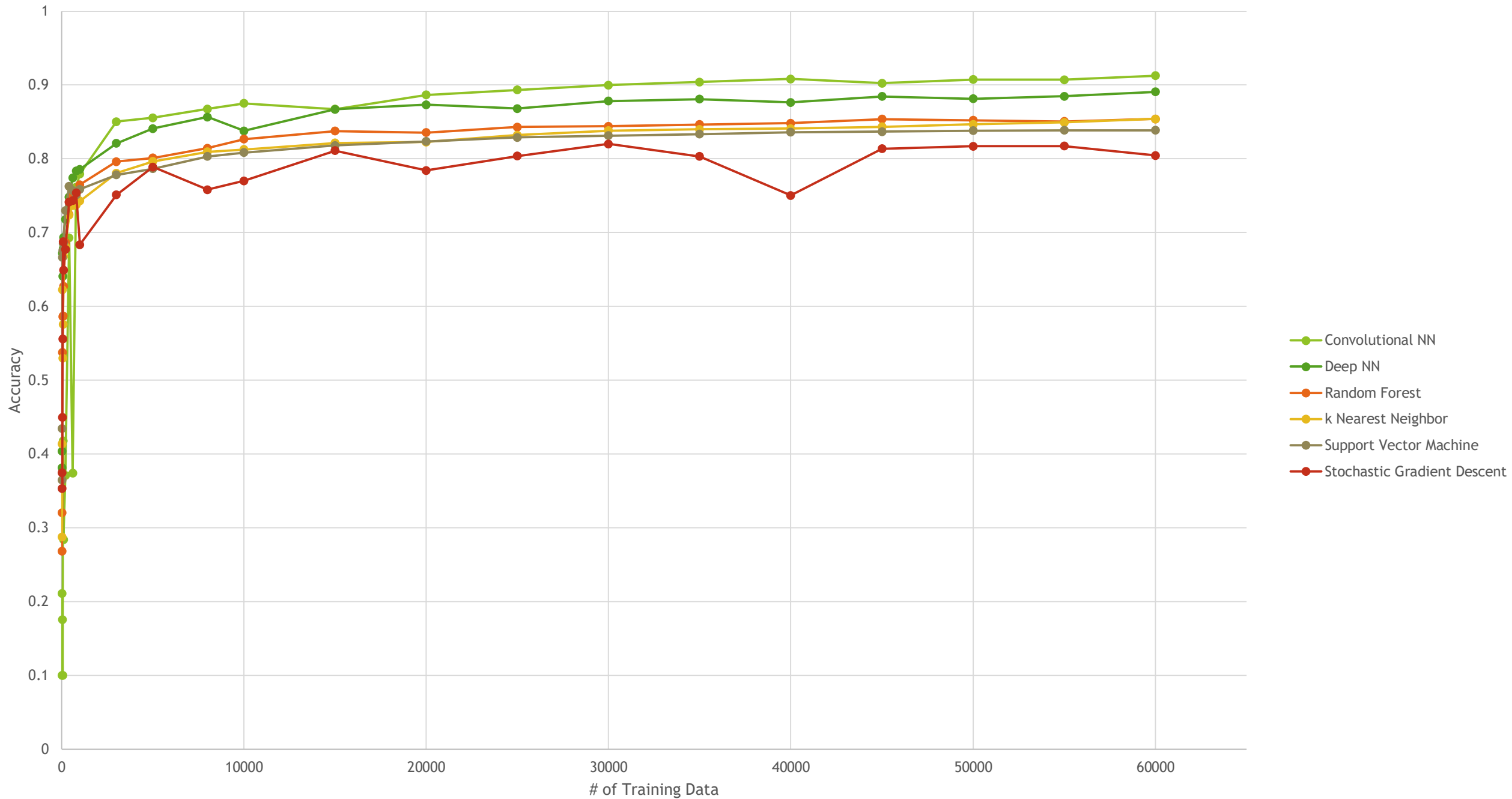
- ▶ Fashion-MNIST Dataset
 - ▶ 70000 labeled 28x28 images - 60000 training, 10000 testing
 - ▶ 10 fashion categories - T-shirt, Dress, Coat, Sandal, etc.
- ▶ MNIST Dataset
 - ▶ 70000 labeled 28x28 images - 60000 training, 10000 testing
 - ▶ 10 handwritten digits
- ▶ Omniglot Dataset
 - ▶ “Transpose of MNIST”
 - ▶ 1623 different handwritten characters from 50 different alphabets - 30 training, 20 testing
 - ▶ 20 examples of every character
 - ▶ Mongolian, Armenian, Hebrew, (Futurama, Alphabet of the Magi)
- ▶ Why these?
 - ▶ Clean - no corrupted files
 - ▶ Image Datasets - intuitive, CNN-applicable
 - ▶ Lots of Documentation, Academic
 - ▶ Good for low-shot learning - Omniglot



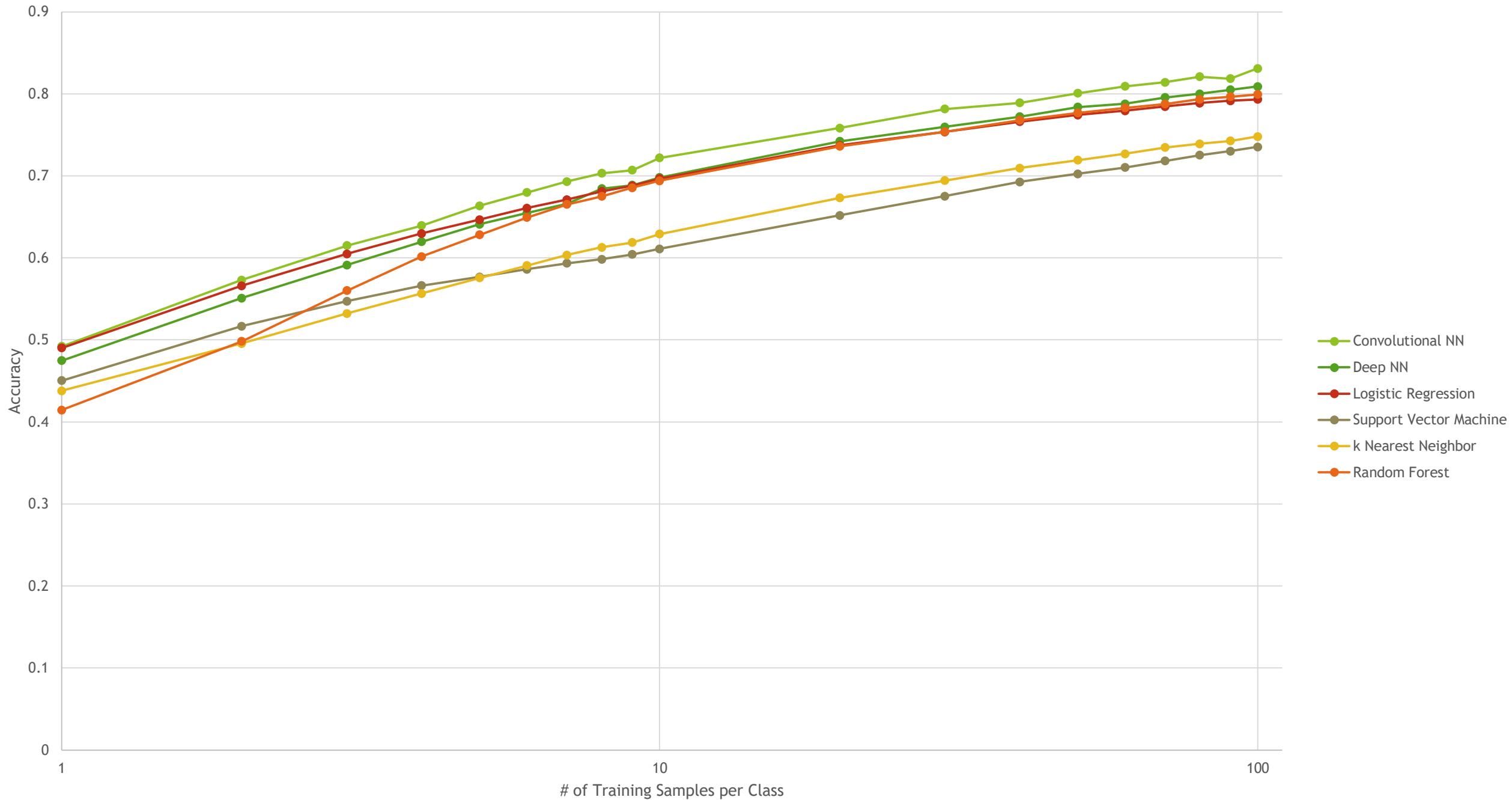
Algorithms tested

- ▶ Convolutional Neural Network
- ▶ Deep Neural Network
- ▶ Random Forest
- ▶ K Nearest-Neighbor
- ▶ Support Vector Machine
- ▶ Logistic Regression
- ▶ ~~Stochastic Gradient Descent~~

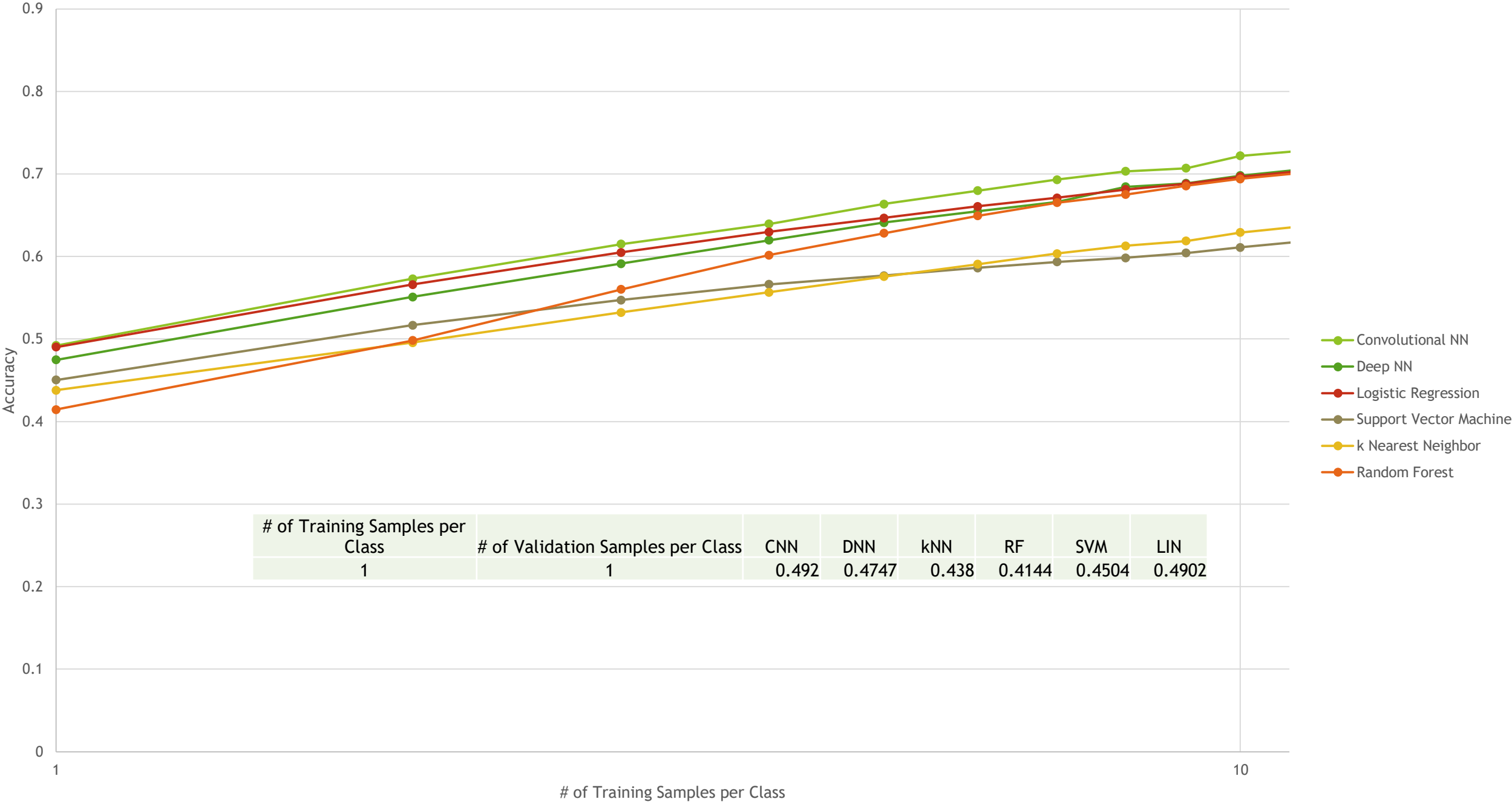
Variable Train Set: 60000 datapoints, 1 trial



Variable Train Set: 100 datapoints, averaged full database trials, optimized, balanced classes



Variable Train Set: 10 datapoints, averaged full database trials, optimized, balanced classes



One-Shot Learning

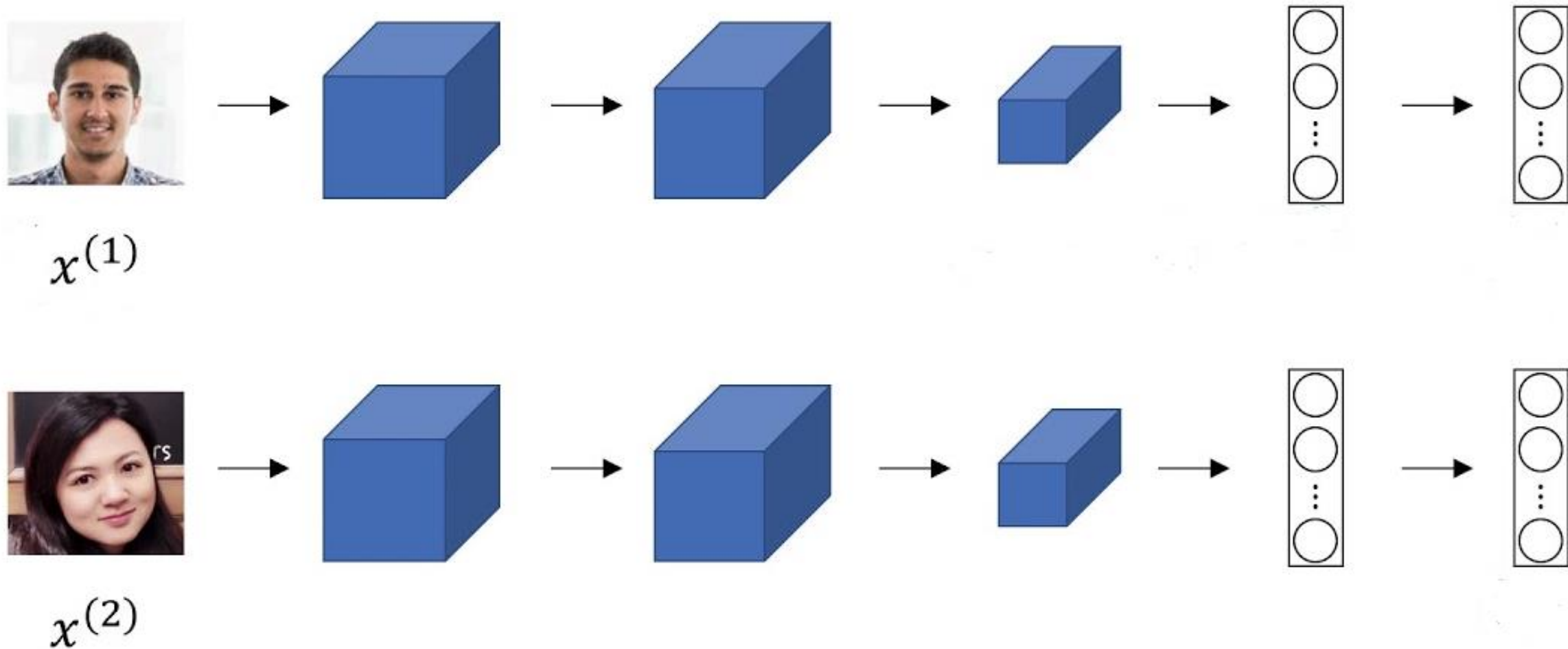


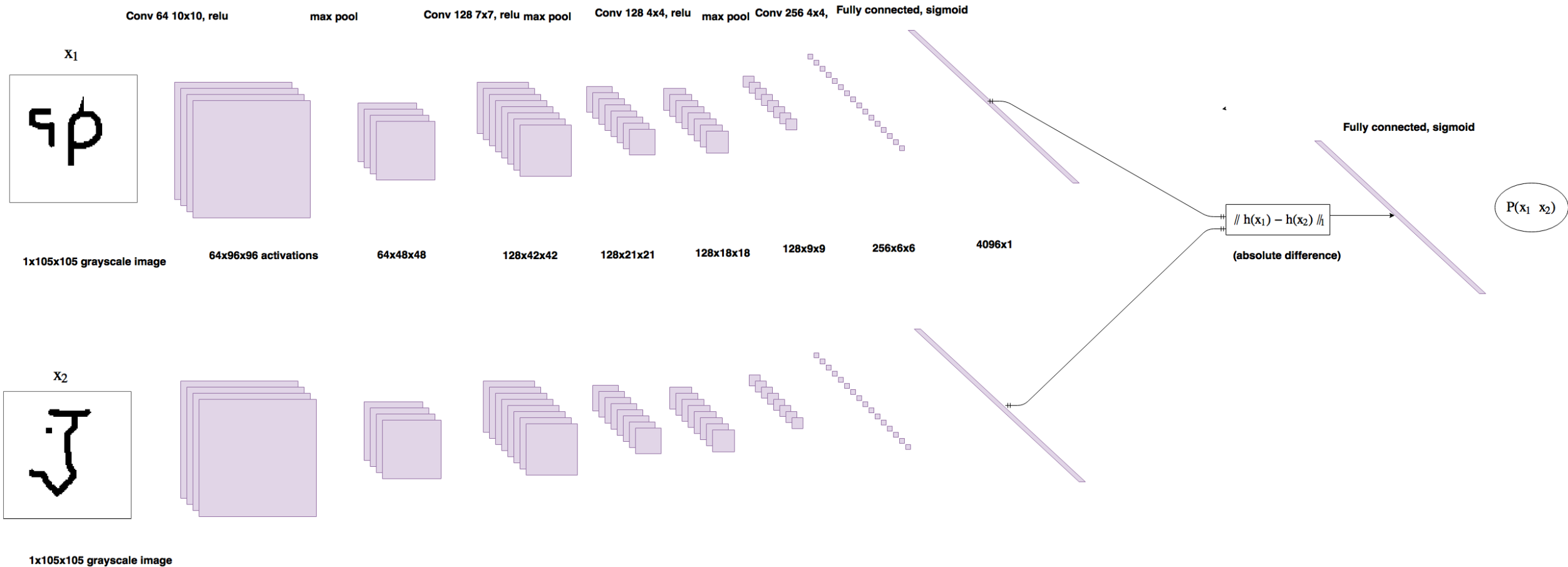
4-way one-shot
learning

Specialized techniques

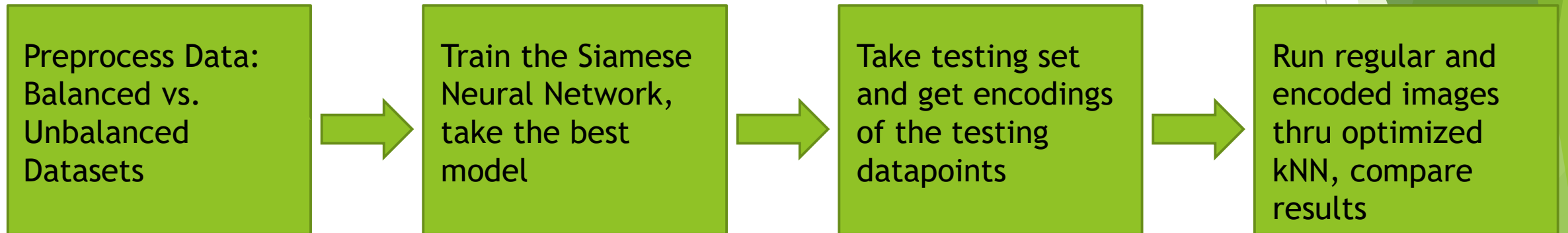
- ▶ **Siamese Network**
- ▶ Triplet Loss
- ▶ CapsuleNet

Siamese Neural Network (SNN)





How we use SNNs



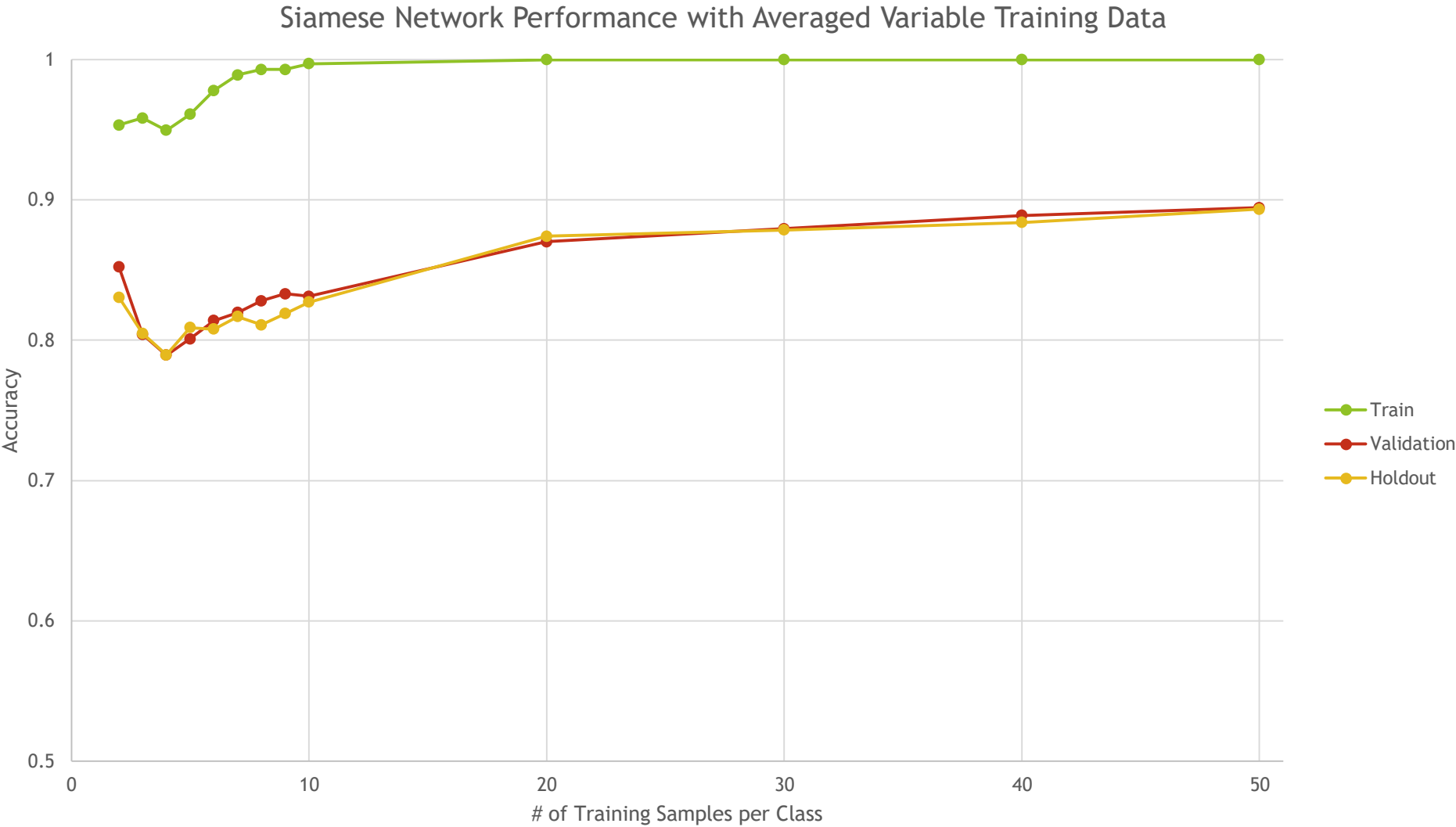
SNN Study 1: Variable Train Set

- ▶ Continuation of the first experiment - how much does the Siamese network improve the accuracy for a variable train set?

SNN Study 1: Variable Train Set - Network Accuracy

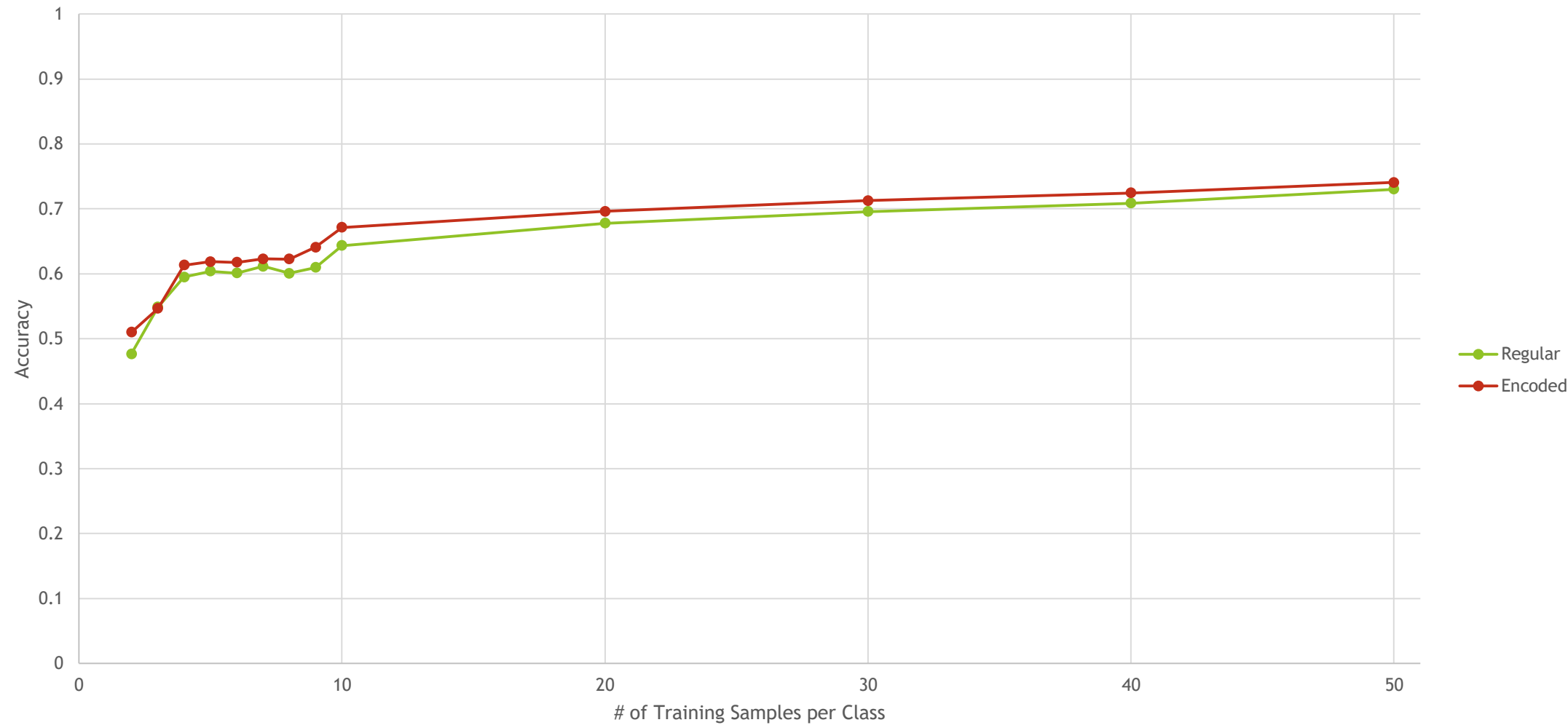
Network Scores:		
Train	Val	Hold
98.67%	86.62%	86.44%

Results for half similar pairs, half different pairs on 50 samples per class



SNN Study 1: Variable Train Set - Testing Accuracy

Averaged Variable Training Data on Siamese Network



SNN Study 2: Using entire dataset

- ▶ 30000 images to train Siamese Network
 - ▶ Subset of 150 images per class
 - ▶ $\binom{1500}{2} = 1124250$ pairs compared to $\binom{500}{2} = 124750$ pairs
- ▶ 30000 images to test accuracies on an optimized kNN

Balanced split

Network Accuracies:			Testing Accuracies:	
Train	Val	Hold	Regular	Encoded
99.99%	89.9%	88.58%	76.96%	77.18%

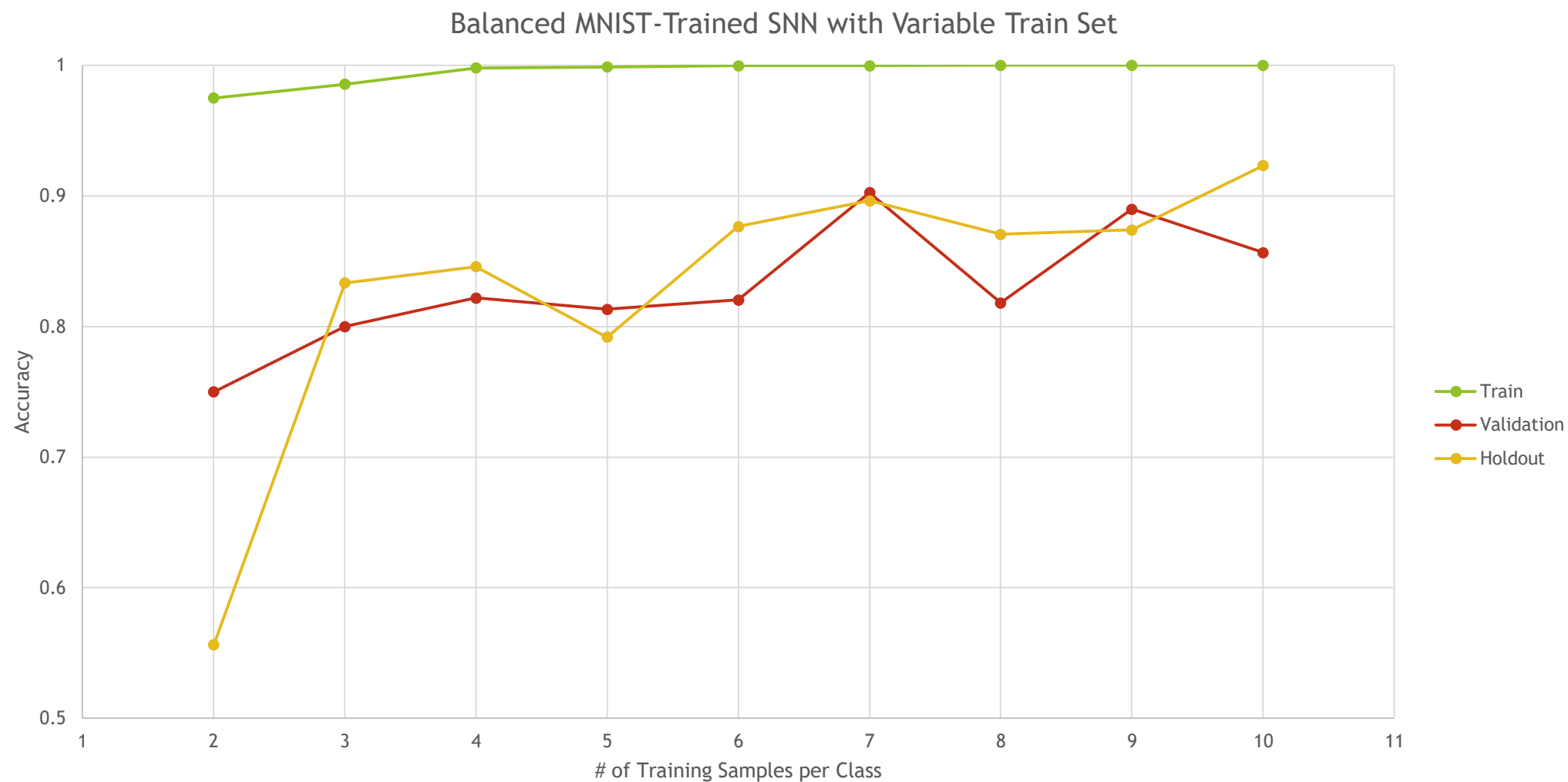
Unbalanced split

Network Accuracies:			Testing Accuracies:	
Train	Val	Hold	Regular	Encoded
99.99%	93.82%	91.73%	76.99%	83.15%

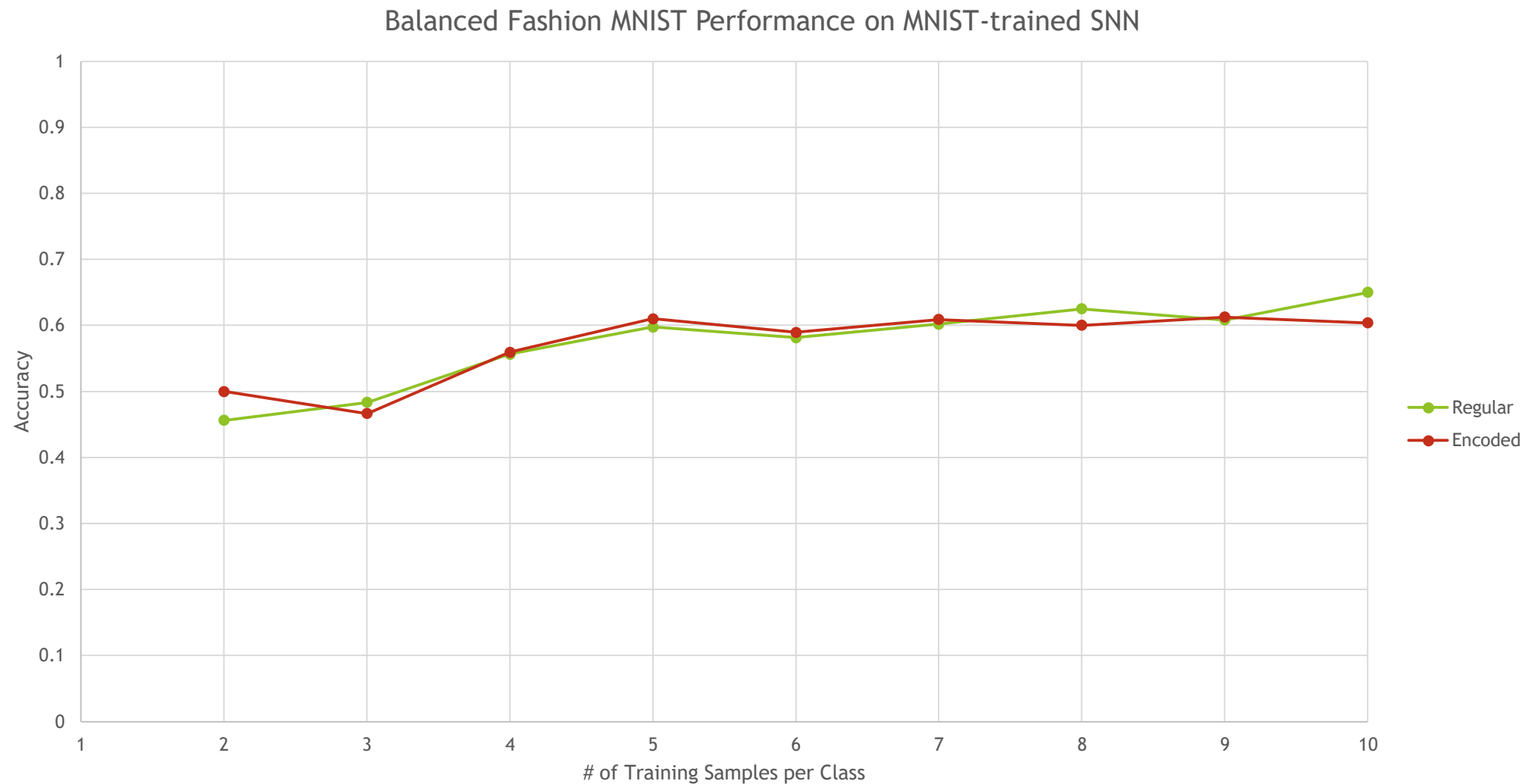
SNN Study 3: Generalization of features

- ▶ Paper: Trained model with omniglot, tested on MNIST - 70.3%
- ▶ MNIST dataset → Fashion MNIST dataset

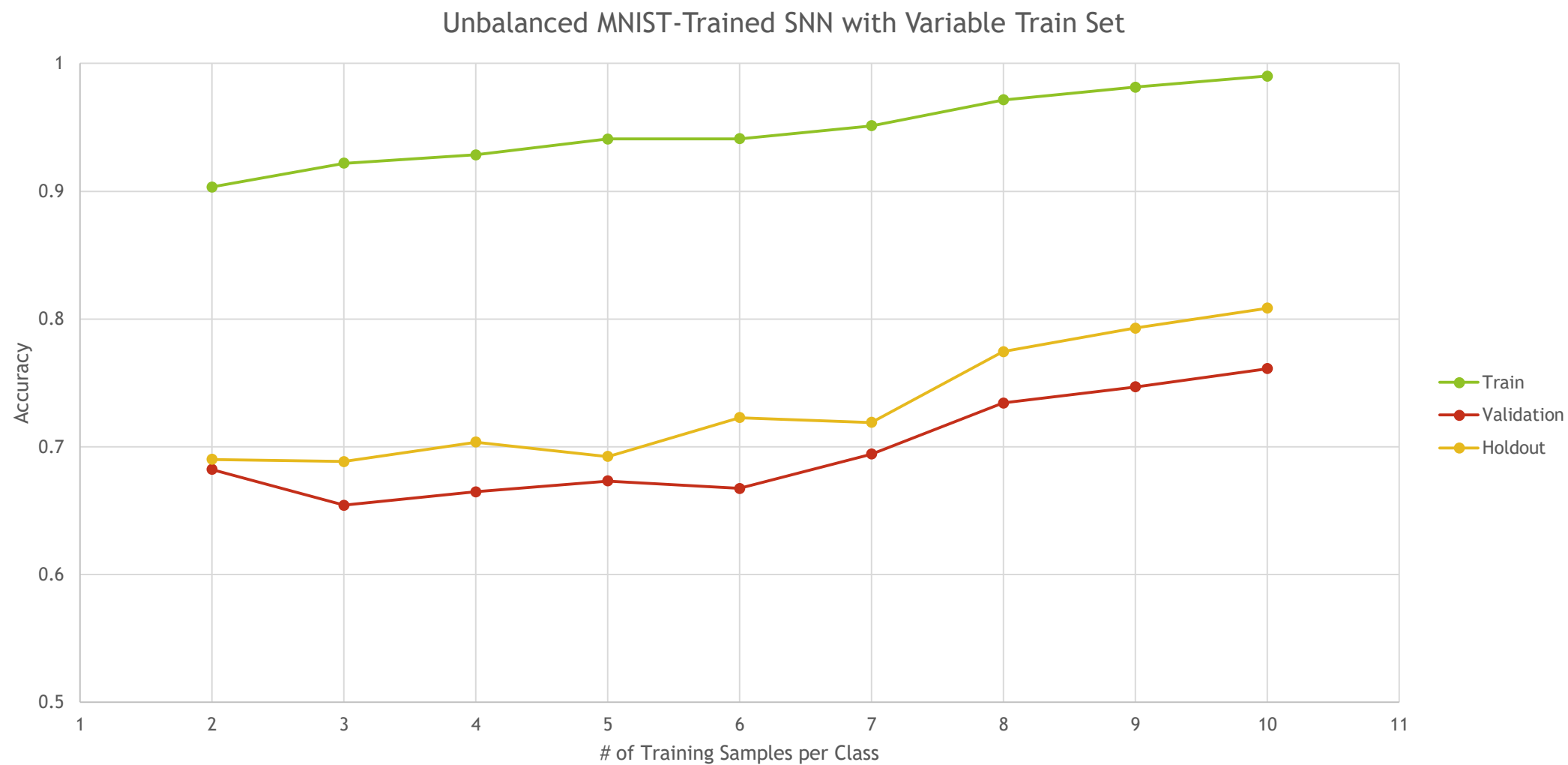
SNN Study 3 - Network Accuracy, Balanced



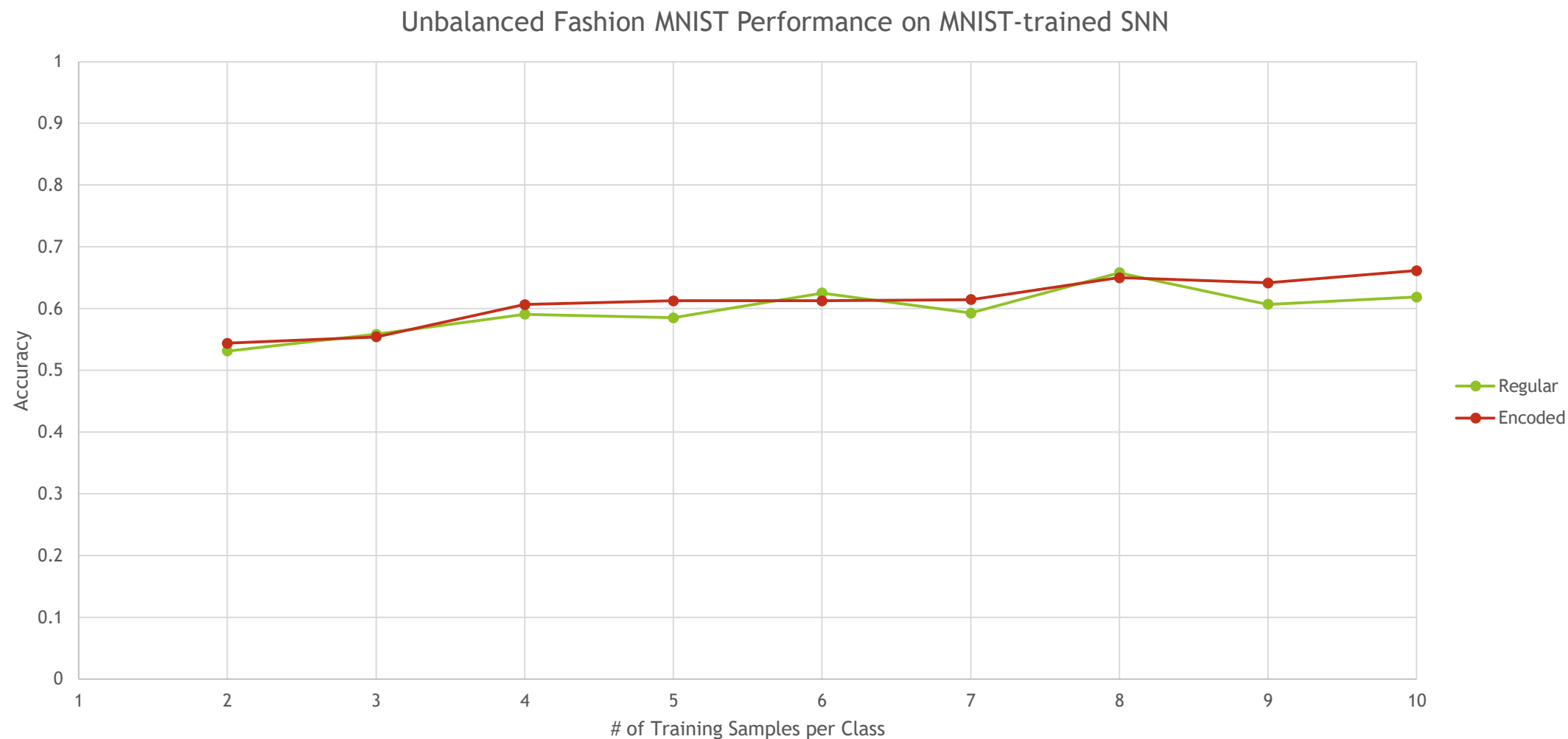
SNN Study 3 - Testing Accuracy, Balanced



SNN Study 3 - Network Accuracy, Unbalanced

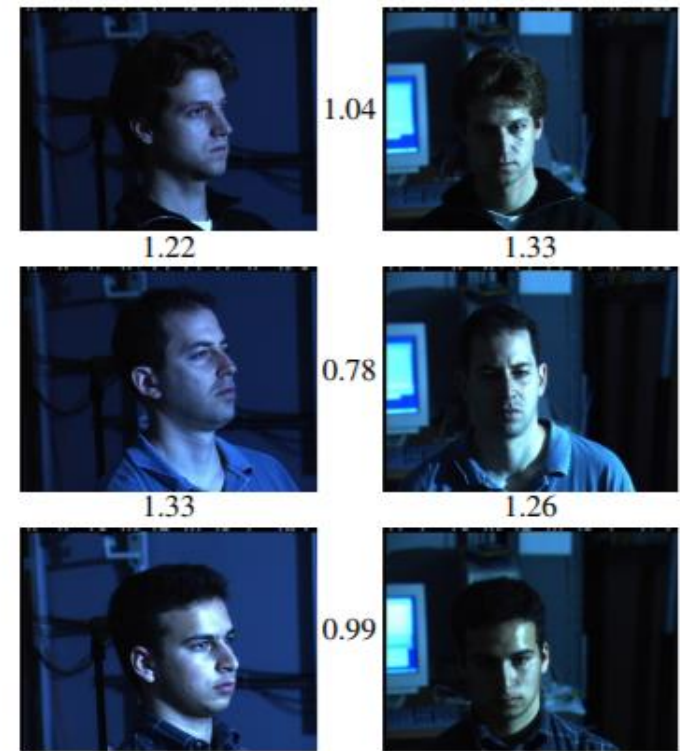


SNN Study 3 - Testing Accuracy, Unbalanced



Future work

- ▶ Another SNN Study: Train with 150 pts per class, test on variable, small dataset
- ▶ Finetune the models for study 3
- ▶ Triplet Loss Algorithm
- ▶ CapsuleNet Algorithm
- ▶ Apply algorithms on Omniglot Dataset
- ▶ Apply algorithms on cancer dataset



Related Papers: One-Shot Learning

- ▶ “Low Data Drug Discovery with One-Shot Learning”
- ▶ “Active One-Shot Learning”

Acknowledgements

Thanks to Argonne National Labs and the Department of Energy for making this project possible. Additional thanks to Dr. Fangfang Xia and Neil Getty for their continuous expertise and support throughout the internship.