

# XCS330 Problem Set 2

---

**Due Sunday, April 16 at 11:59pm PT.**

## Guidelines

1. If you have a question about this homework, we encourage you to post your question on our Slack channel, at <http://xcs330-scpd.slack.com/>
2. Familiarize yourself with the collaboration and honor code policy before starting work.
3. For the coding problems, you must use the packages specified in the provided environment description. Since the autograder uses this environment, we will not be able to grade any submissions which import unexpected libraries.

## Submission Instructions

**Written Submission:** Some questions in this assignment require a written response. For these questions, you should submit a PDF with your solutions online in the online student portal. As long as the PDF is legible and organized, the course staff has no preference between a handwritten and a typeset L<sup>A</sup>T<sub>E</sub>X submission. If you wish to typeset your submission and are new to L<sup>A</sup>T<sub>E</sub>X, you can get started with the following:

- Type responses only in `submission.tex`.
- Submit the compiled PDF, **not** `submission.tex`.
- Use the commented instructions within the `Makefile` and `README.md` to get started.

**Coding Submission:** Some questions in this assignment require a coding response. For these questions, you should submit only the `src/submission.py` file in the online student portal. For further details, see Writing Code and Running the Autograder below.

## Honor code

We strongly encourage students to form study groups. Students may discuss and work on homework problems in groups. However, each student must write down the solutions independently, and without referring to written notes from the joint session. In other words, each student must understand the solution well enough in order to reconstruct it by him/herself. In addition, each student should write on the problem set the set of people with whom s/he collaborated. Further, because we occasionally reuse problem set questions from previous years, we expect students not to copy, refer to, or look at the solutions in preparing their answers. It is an honor code violation to intentionally refer to a previous year's solutions. More information regarding the Stanford honor code can be found at <https://communitystandards.stanford.edu/policies-and-guidance/honor-code>.

## Writing Code and Running the Autograder

All your code should be entered into `src/submission.py`. When editing `src/submission.py`, please only make changes between the lines containing `### START_CODE_HERE ###` and `### END_CODE_HERE ###`. Do not make changes to files other than `src/submission.py`.

The unit tests in `src/grader.py` (the autograder) will be used to verify a correct submission. Run the autograder locally using the following terminal command within the `src/` subdirectory:

```
$ python grader.py
```

There are two types of unit tests used by the autograder:

- **basic:** These tests are provided to make sure that your inputs and outputs are on the right track, and that the hidden evaluation tests will be able to execute.

- **hidden:** These unit tests are the evaluated elements of the assignment, and run your code with more complex inputs and corner cases. Just because your code passed the basic local tests does not necessarily mean that they will pass all of the hidden tests. These evaluative hidden tests will be run when you submit your code to the Gradescope autograder via the online student portal, and will provide feedback on how many points you have earned.

For debugging purposes, you can run a single unit test locally. For example, you can run the test case `3a-0-basic` using the following terminal command within the `src/` subdirectory:

```
$ python grader.py 3a-0-basic
```

Before beginning this course, please walk through the [Anaconda Setup for XCS Courses](#) to familiarize yourself with the coding environment. Use the env defined in `src/environment.yml` to run your code. This is the same environment used by the online autograder.

## Test Cases

The autograder is a thin wrapper over the python `unittest` framework. It can be run either locally (on your computer) or remotely (on SCPD servers). The following description demonstrates what test results will look like for both local and remote execution. For the sake of example, we will consider two generic tests: `1a-0-basic` and `1a-1-hidden`.

### Local Execution - Hidden Tests

All hidden tests rely on files that are not provided to students. Therefore, the tests can only be run remotely. When a hidden test like `1a-1-hidden` is executed locally, it will produce the following result:

```
----- START 1a-1-hidden: Test multiple instances of the same word in a sentence.
----- END 1a-1-hidden [took 0:00:00.011989 (max allowed 1 seconds), ???/3 points] (hidden test ungraded)
```

### Local Execution - Basic Tests

When a basic test like `1a-0-basic` passes locally, the autograder will indicate success:

```
----- START 1a-0-basic: Basic test case.
----- END 1a-0-basic [took 0:00:00.000062 (max allowed 1 seconds), 2/2 points]
```

When a basic test like `1a-0-basic` fails locally, the error is printed to the terminal, along with a stack trace indicating where the error occurred:

```
----- START 1a-0-basic: Basic test case.
<class 'AssertionError'>
{'a': 2, 'b': 1} != None ← This error caused the test to fail.
File "/Users/grinch/Local_Documents/Software/anaconda3/envs/XCS221/lib/python3.6/unittest/case.py", line 59, in testPartExecutor
yield
File "/Users/grinch/Local_Documents/Software/anaconda3/envs/XCS221/lib/python3.6/unittest/case.py", line 605, in run
testMethod()
File "/Users/grinch/Local_Documents/SCPD/XCS221/A1/src/graderUtil.py", line 54, in wrapper
result = func(*args, **kwargs)
File "/Users/grinch/Local_Documents/SCPD/XCS221/A1/src/graderUtil.py", line 83, in wrapper
result = func(*args, **kwargs)
File "/Users/grinch/Local_Documents/SCPD/XCS221/A1/src/grader.py", line 23, in test_0
submission.extractWordFeatures("a b a") ← In this case, start your debugging
in line 23 of grader.py.
File "/Users/grinch/Local_Documents/Software/anaconda3/envs/XCS221/lib/python3.6/unittest/case.py", line 829, in assertEqual
assertion_func(first, second, msg=msg)
File "/Users/grinch/Local_Documents/Software/anaconda3/envs/XCS221/lib/python3.6/unittest/case.py", line 822, in _baseAssertEqual
raise self.failureException(msg)
----- END 1a-0-basic [took 0:00:00.003809 (max allowed 1 seconds), 0/2 points]
```

## Remote Execution

Basic and hidden tests are treated the same by the remote autograder. Here are screenshots of failed basic and hidden tests. Notice that the same information (error and stack trace) is provided as the in local autograder, now for both basic and hidden tests.

## 1a-0-basic) Basic test case. (0.0/2.0)

```
<class 'AssertionError': {'a': 2, 'b': 1} != None
File "/autograder/source/miniconda/envs/XCS221/lib/python3.6/unittest/case.py", line 59, in testPartExecutor
    yield
File "/autograder/source/miniconda/envs/XCS221/lib/python3.6/unittest/case.py", line 605, in run
    testMethod()
File "/autograder/source/graderUtil.py", line 54, in wrapper
    result = func(*args, **kwargs)
File "/autograder/source/graderUtil.py", line 83, in wrapper
    result = func(*args, **kwargs)
File "/autograder/source/grader.py", line 23, in test_0
    submission.extractWordFeatures("a b a"))
File "/autograder/source/miniconda/envs/XCS221/lib/python3.6/unittest/case.py", line 829, in assertEqual
    assertion_func(first, second, msg=msg)
File "/autograder/source/miniconda/envs/XCS221/lib/python3.6/unittest/case.py", line 822, in _baseAssertEqual
    raise self.failureException(msg)
```

Just like in the local autograder, this error caused the test to fail.

Just like in the local autograder, start your debugging in line 23 of grader.py.

## 1a-1-hidden) Test multiple instances of the same word in a sentence. (0.0/3.0)

```
<class 'AssertionError': {'a': 23, 'ab': 22, 'aa': 24, 'c': 16, 'b': 15} != None
File "/autograder/source/miniconda/envs/XCS221/lib/python3.6/unittest/case.py", line 59, in testPartExecutor
    yield
File "/autograder/source/miniconda/envs/XCS221/lib/python3.6/unittest/case.py", line 605, in run
    testMethod()
File "/autograder/source/graderUtil.py", line 54, in wrapper
    result = func(*args, **kwargs)
File "/autograder/source/graderUtil.py", line 83, in wrapper
    result = func(*args, **kwargs)
File "/autograder/source/grader.py", line 31, in test_1
    self.compare_with_solution_or_wait(submission, 'extractWordFeatures', lambda f: f(sentence))
File "/autograder/source/graderUtil.py", line 183, in compare_with_solution_or_wait
    self.assertEqual(ans1, ans2)
File "/autograder/source/miniconda/envs/XCS221/lib/python3.6/unittest/case.py", line 829, in assertEqual
    assertion_func(first, second, msg=msg)
File "/autograder/source/miniconda/envs/XCS221/lib/python3.6/unittest/case.py", line 822, in _baseAssertEqual
    raise self.failureException(msg)
```

This error caused the test to fail.

Start your debugging in line 31 of grader.py.

Finally, here is what it looks like when basic and hidden tests pass in the remote autograder.

## 1a-0-basic) Basic test case. (2.0/2.0)

## 1a-1-hidden) Test multiple instances of the same word in a sentence. (3.0/3.0)

## Data Processing and Black-Box Meta-Learning

**Goals:** In this assignment, we will look at meta-learning for few shot classification. You will:

1. Learn how to process and partition data for meta learning problems, where training is done over a distribution of training tasks  $p(\mathcal{T})$ .
2. Implement and train memory augmented neural networks, a black-box meta-learner that uses a recurrent neural network [1].
3. Analyze the learning performance for different size problems.
4. Experiment with model parameters and explore how they improve performance.

We will be working with Omniglot [2], a dataset with 1623 characters from 50 different languages. Each character has 20 28x28 images. We are interested in training models for  $K$ -shot,  $N$ -way classification, i.e. training a classifier to distinguish between  $N$  previously unseen characters, given only  $K$  labeled examples of each character.

For this assignment, please submit the following files to gradescope to receive points for coding questions:

- `submission.pdf`
- `src/submission/__init__.py`
- `src/submission/load_data.py`
- `src/submission/mann.py`
- `src/mann_results_1_2.npy`
- `src/mann_results_2_2.npy`
- `src/mann_results_1_3.npy`
- `src/mann_results_1_4.npy`

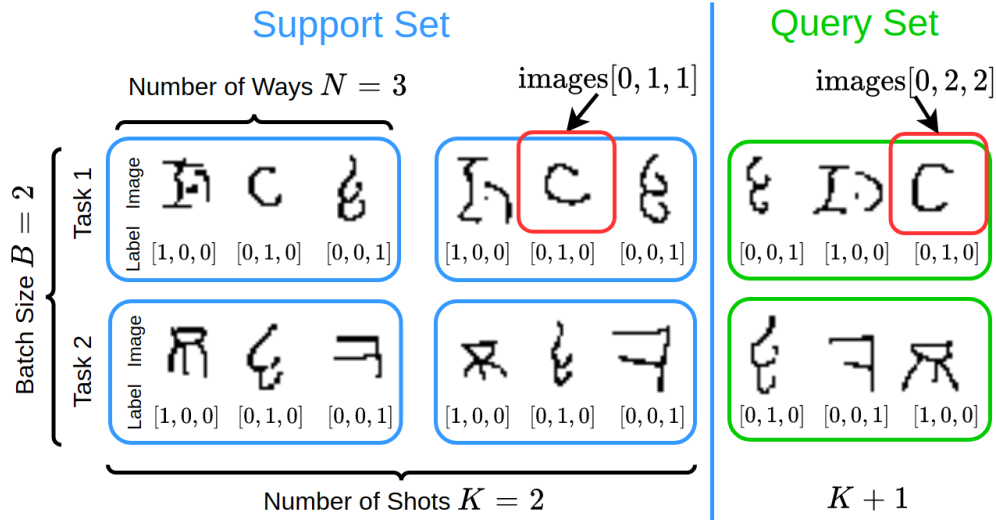


Figure 1: Example data batch from the Data Generator. The first  $K$  sets of images form the support set and are passed in the *same order*. The final set of images forms the query set and must be shuffled.

### 1. [15 points (Coding)] Data Processing for Few-Shot Classification

Before training any models, you must write code to sample batches for training. Fill in the `_sample` function in the `DataGenerator` class. The class already has variables defined for batch size `batch_size` ( $B$ ), number of classes `num_classes` ( $N$ ), and number of samples per class `num_samples_per_class` ( $K + 1$ ). Your code should:

- Sample  $N$  different characters from either the specified train, test, or validation folder.
- Load  $K + 1$  images per character and collect the associated labels, using  $K$  images per class for the support set and 1 image per class for the query set.
- Format the data and return two tensors, one of flattened images with shape  $[K + 1, N, 784]$  and one of one-hot labels  $[K + 1, N, N]$ .

Note that your code only needs to return one single (image, label) tuple. We batch the inputs using an instance of `torch.utils.data.DataLoader`, and the final shape input images is  $[B, K + 1, N, 784]$ , and that of the input labels is  $[B, K + 1, N, N]$ , where  $B$  is the batch size.

Figure 1 illustrates the data organization. In this example, we have: (1) images from  $N = 3$  different classes; (2) we are provided  $K = 2$  sets of labeled images in the support set and (3) our batch consists of only two tasks, i.e.  $B = 2$ .

- We will sample both the support and query sets as a single batch, hence one batch element should obtain image and label tensors of shapes  $[K + 1, N, 784]$  and  $[K + 1, N, N]$  respectively. In the example of Fig. 1, `images[0, 1, 1]` would be the image of the letter "C" in the support set with corresponding class label `[0, 1, 0]` and `images[0, 2, 2]` would be the the letter "C" in the query set (with the same label).
- We must shuffle the order of examples in the **query set**, as otherwise the network can learn to output the same sequence of classes and achieve 100% accuracy, without actually learning to recognize the images. If you get 100% accuracy, you likely did not shuffle the query data correctly. In principle, you should be able to shuffle the order of data in the support set as well; however, this makes the model optimization much harder. **You should feed the support set examples in the same, fixed order**. In the example above, the support set examples are always in the same order.

We provide helper functions to (1) take a list of folders and provide paths to image files/labels, and (2) to take an image file path and return a flattened numpy matrix. The functions `np.random.shuffle` and `np.eye` will also be helpful. **Be careful about output shapes and data types!**

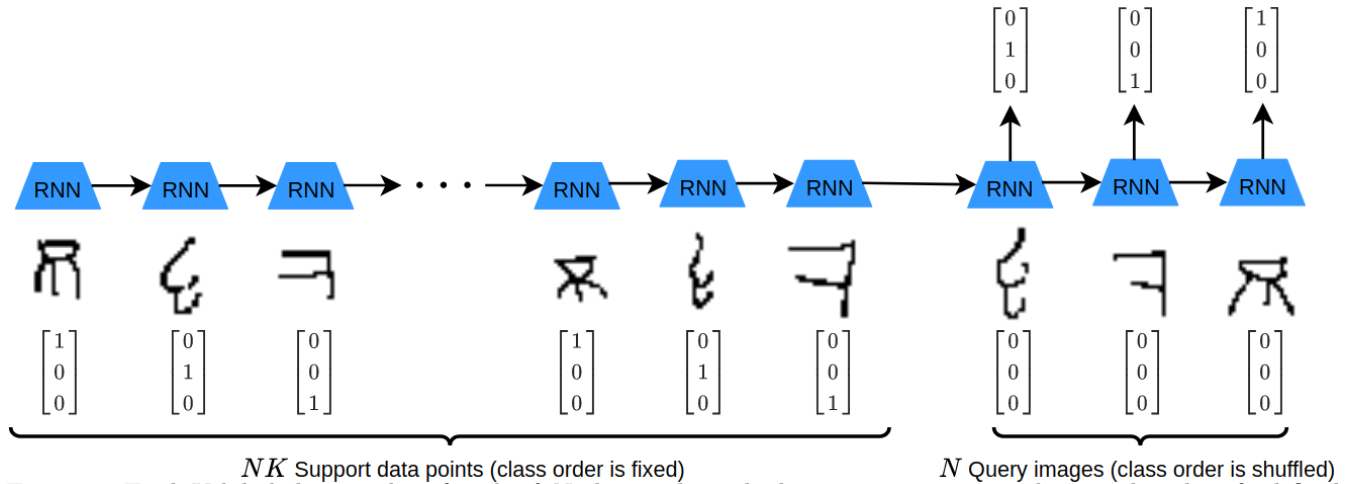


Figure 2: Feed  $K$  labeled examples of each of  $N$  classes through the memory-augmented network. Then feed final set of  $N$  examples and optimize to minimize loss.

## 2. [15 points (Coding)] Memory Augmented Neural Networks (MANN) [1, 3]

We will now be implementing few-shot classification using memory augmented neural networks (MANNs). The main idea of MANN is that the network should learn how to encode the first  $K$  examples of each class into memory such that it can be used to accurately classify the  $K + 1$ th example. See Figure 2 for a graphical representation of this process.

Data processing will be done as in SNAIL [3]. Each set of labels and images are concatenated together, and the  $N * K$  support set examples are sequentially passed through the network as shown in Fig. 2. Then the query example of each class is fed through the network, **concatenated with 0 instead of the true label**. The loss is computed between the query set predictions and the ground truth labels, which is then backpropagated through the network. **Note:** The loss is *only* computed on the set of  $N$  query images, which comprise of the last examples from each character class.

In the `submission/mann.py` file:

- Fill in the `forward` function of the `MANN` class to take in image tensor of shape  $[B, K + 1, N, 784]$  and a label tensor of shape  $[B, K + 1, N, N]$  and output labels of shape  $[B, K + 1, N, N]$ . The layers to use have already been defined for you in the `__init__` function. *Hint: Remember to pass zeros, not the ground truth labels for the final  $N$  examples.*
- Fill in the function called `loss_function` in the `MANN` class which takes as input the  $[B, K + 1, N, N]$  labels and  $[B, K + 1, N, N]$  predicted labels and computes [the cross entropy loss](#) only on the  $N$  test images.

**Note:** Both of the above functions will need to be backpropagated through, so they need to be written in PyTorch in a differentiable way.

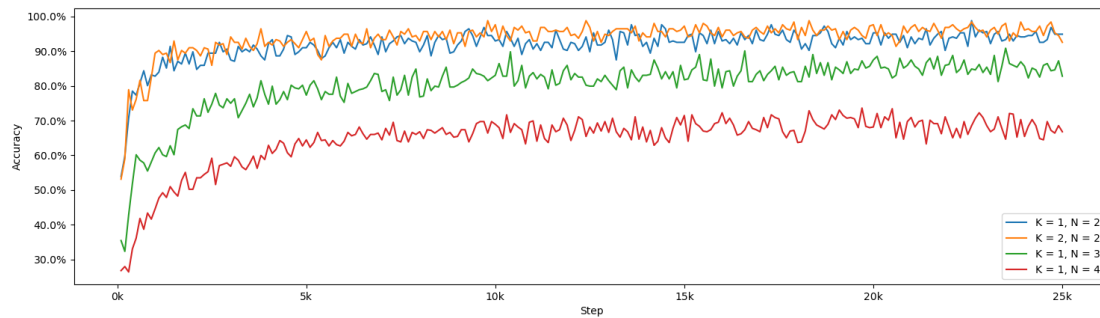


Figure 3: Tensorboard plot of the meta-test query set classification accuracy over training iterations for all configurations.

### 3. [12 points (Coding)] [4 points (Written)] Train and Analysis

Once you have completed problems 1 and 2, you can train your few shot classification model. You should observe both the support and query losses go down, and the query accuracy go up. Now we will examine how the performance varies for different size problems. Train models for the following values of  $K$  and  $N$ :

- $K = 1, N = 2$
- $K = 2, N = 2$
- $K = 1, N = 3$
- $K = 1, N = 4$

Example code:

```
python main.py --num_shot K --num_classes N
```

For checking training results and/or taking a screenshot for the writeup, use:

```
tensorboard --logdir runs/
```

You should start with the case  $K = 1, N = 2$  as it can aid you in the implementation and debugging process. Your model should be able to achieve a query set accuracy of above 90% in this first two scenarios scenario on held-out test tasks, around 80% in the third scenario, and around 70% in the final scenario.

We have provided in Figure 3 a unified view of the meta-test query set classification accuracy over training iterations for all above mentioned configurations. Answer the following questions:

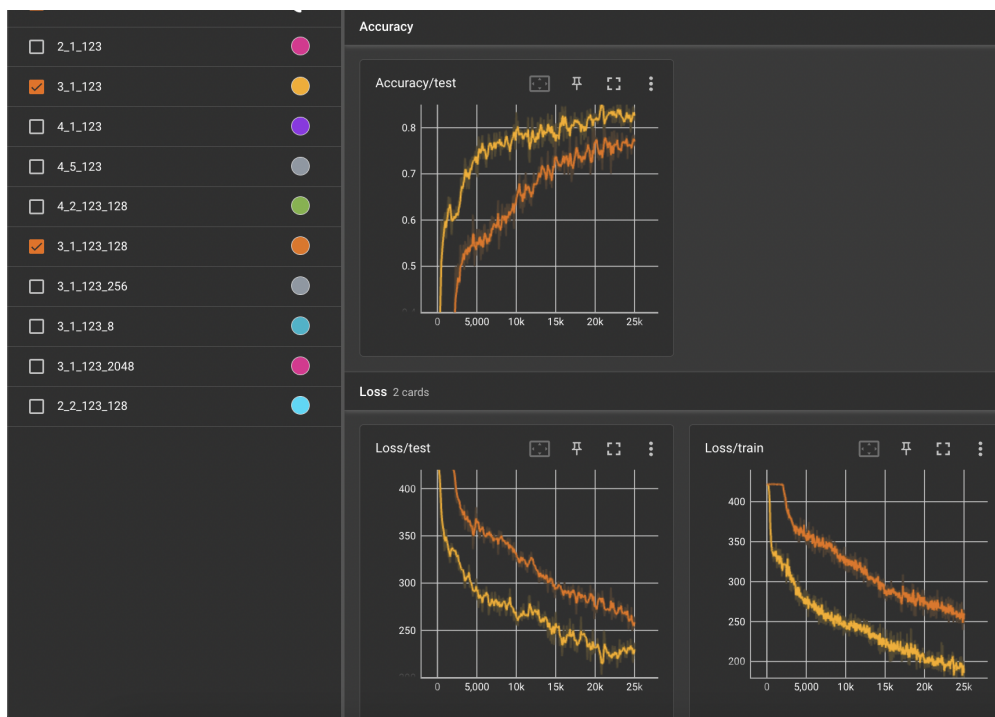
- How does increasing the number of classes affect learning and performance?
- How does increasing the number of examples in the support set affect performance?



## 4. [6 points (Written, Extra Credit)] Experimentation

- a Experiment with one hyperparameter that affects the performance of the model, such as the type of recurrent layer, size of hidden state, learning rate, or number of layers.

The plot that shows how the meta-test query set classification accuracy of the model changes on 1-shot, 3-way classification as you change the parameter should look as follows.



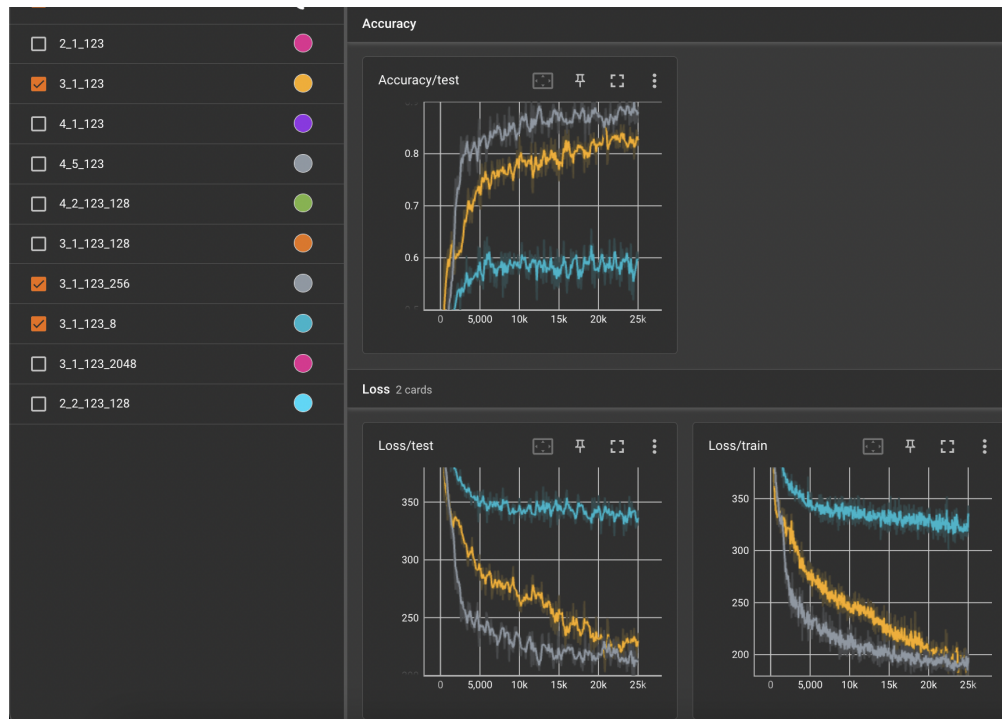
Provide a brief rationale for why you chose the parameter and what you observed in the caption for the plot.

- b **Extra Credit:** In this question we'll explore the effect of memory representation on model performance. We will focus on the  $K = 1$ ,  $N = 3$  case.

In the previous experiments we used an LSTM model with 128 units. Consider additional memory sizes of 256, and 8.

The plots should look as follows:





How does increasing and decreasing the memory capacity influence performance?

## Coding Deliverables

## References

- [1] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1842–1850, New York, New York, USA, 20–22 Jun 2016. PMLR.
- [2] Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- [3] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. Meta-learning with temporal convolutions. *CoRR*, abs/1707.03141, 2017.

This handout includes space for every question that requires a written response. Please feel free to use it to handwrite your solutions (legibly, please). If you choose to typeset your solutions, the `README.md` for this assignment includes instructions to regenerate this handout with your typeset L<sup>A</sup>T<sub>E</sub>X solutions.

---

3.a

3.b

4.a

4.b