# Syllabus for CS6787

## Advanced Machine Learning Systems — Spring 2024

| | | | |
|---|---|---|---|
| **Term** | Spring 2024 | **Instructor** | Christopher De Sa |
| **Room** | Phillips Hall 101 | **E-mail** | [email hidden] |
| **Schedule** | MW 7:30pm – 8:45pm | **Office hours** | W 2:30pm – 3:30pm |
| **Forum** | Ed Discussion | **Office** | Gates 426 |

So you've taken a machine learning class. You know the models people use to solve their problems. You know the algorithms they use for learning. You know how to evaluate the quality of their solutions.

But when we look at a large-scale machine learning application that is deployed in practice, it's not always exactly what you learned in class. Sure, the basic models, the basic algorithms are all there. But they're modified a bit, in a bunch of different ways, to run faster and more efficiently. And these modifications are really important—they often are what make the system tractable to run on the data it needs to process.

CS6787 is a graduate-level introduction to these system-focused aspects of machine learning, covering guiding principles and commonly used techniques for scaling up learning to large data sets. Informally, we will cover the techniques that lie between a standard machine learning course and an efficient systems implementation: both statistical/optimization techniques based on improving the convergence rate of learning algorithms and techniques that improve performance by leveraging the capabilities of the underlying hardware. Topics will include stochastic gradient descent, acceleration, variance reduction, methods for choosing hyperparameters, parallelization within a chip and across a cluster, popular ML frameworks, and innovations in hardware architectures. An open-ended project in which students apply these techniques is a major part of the course.

**Prerequisites:** Knowledge of machine learning at the level of CS4780. *If you are an undergraduate, you should have taken CS4780 or an equivalent course, since it is a prerequisite.* Knowledge of computer systems and hardware on the level of CS 3410 is recommended, but this is not a prerequisite.

**Format:** About half of the classes will involve traditionally formatted lectures. For the other half of the classes, we will read and discuss two seminal papers relevant to the course topic. These classes will involve presentations by groups of students of the paper contents (each student will sign up in a group to present one paper for 15-20 minutes) followed by breakout discussions about the material. Historically, the lectures have occurred on Mondays and the discussions have occurred on Wednesdays, but due to the non-standard timeline this semester, these course elements will be scheduled irregularly (see schedule below).

**Grading:** Students will be evaluated on the following basis.

| | |
|---|---|
| 20% | Paper presentation |
| 10% | Discussion participation |
| 20% | Paper reviews |
| 10% | Programming assignments |
| 40% | Final project |

**Paper review parameters:** Paper reviews should be about one page (single-spaced) in length. The review guidelines should mirror what an actual conference review would look like (although you needn't assign scores or anything like that). In particular you should at least: (1) summarize the paper, (2) discuss the paper's strengths and weaknesses, and (3) discuss the paper's impact. For reference, you can read the ICML reviewer guidelines. Of course, your review will not be precisely like a real review, in large part because we already know the impact of these papers. You can submit any review up to two days late with no penalty. Students who presented a paper do not have to submit a review of that paper (although you can if you want).

**Final project parameters (subject to change):** The final project can be done in groups of up to three (although more work will be expected from groups with more people). The subject of the project is open-ended, but it must include:

- the **implementation of a machine learning system** for some task,
- exploring one or more of the **techniques discussed in the course** (or similar techniques subject to instructor approval),
- to **empirically evaluate the performance** and compare it with some baseline method, in two ways:
  - statistical performance (e.g. iterations to converge to some accuracy threshold), and
  - hardware performance (e.g. throughput or wall-clock time).

The project proposal should satisfy the following constraints:
- The main body should be about one page in length.
- It should describe the project you intend to do.
- It should contain at least one citation of a relevant paper that we did not cover in class (but preferably more).

- It should include some preliminary or exploratory work you've already done, that helps to support the idea that your project is feasible (this preliminary work can be very minimal, but should indicate that you've got started—or at least have a clear idea how to do so).
- In addition to the one-page text proposal, it should contain one short **experiment plan** per person, which should consist of:
  - a hypothesis
  - a proxy statement which describes what metric you are going to use to measure the variables you care about
  - a short protocol statement describing what you are going to do
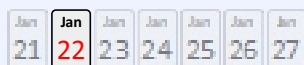  - the results you expect to get

The experiment plan should not be longer than half a page, and may be much shorter.

The project will culminate in a project report of at least four pages, not including references. The project report should be formatted similarly to a workshop paper, and should use the ICML 2019 style or a similar style. The project proposal is due on **Monday, March 25, 2024**. A draft of the final abstract is due for presentation and discussion in class on **Monday, April 29, 2024**. Per the registrar, the final project report is due on **May 15, 2024 at 4:30 PM**.

---

# Course Calendar

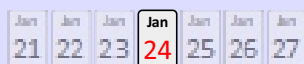| | |
|---|---|
| Monday, January 22<br>**In Person**<br>Jan 21  Jan **22**  Jan 23  Jan 24  Jan 25  Jan 26  Jan 27 | **Lecture #1: Overview.**<br>[Slides] [Demo Notebook]<br>• Overview<br>• Course outline and syllabus<br>• Learning with gradient descent<br>• Stochastic gradient descent: the workhorse of machine learning<br>• Theory of SGD for convex objectives: our first look at trade-offs |
| Wednesday, January 24<br>**In Person**<br>Jan 21  Jan 22  Jan 23  Jan **24**  Jan 25  Jan 26  Jan 27 | **Lecture #2: Backpropagation & ML Frameworks.**<br>[Slides] [Demo Notebook]<br>• Backpropagation and automatic differentiation<br>• Machine learning frameworks I: the user interface<br>• Overfitting<br>• Generalization error<br>• Early stopping<br><br>Optional extra reading. Some older papers on SGD and backpropagation!<br>• Hinton, Geoffrey E. Learning distributed representations of concepts. Proceedings of the eighth annual conference of the cognitive science society. Vol. 1. 1986.<br>• Rumelhart, David E., Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. Cognitive modeling 5.3 (1988): 1.<br>• Tong Zheng. Solving large scale linear prediction problems using stochastic gradient descent algorithms. *Proceedings of the International Conference on Machine Learning (ICML)*, 2004. |
| Monday, January 29<br>**In Person**<br>Jan 28  Jan **29**  Jan 30  Jan 31  Feb 1  Feb 2  Feb 3 | **Lecture #3: Hyperparameters and Tradeoffs.**<br>[Slides] [Demo Notebook]<br>• Our first hyperparameters: step size/learning rate, minibatch size<br>• Regularization<br>• Application-specific forms of regularization<br>• The condition number<br>• Momentum and acceleration<br>• Momentum for quadratic optimization<br>• Momentum for convex optimization<br><br>**Released: Programming Assignment 1.** |
| Wednesday, January 31<br>**In Person**<br>Jan 28  Jan 29  Jan 30  Jan **31**  Feb 1  Feb 2  Feb 3 | **Paper Discussion 1a.**<br>Attention is all you need<br>Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin.<br>*In Advances in neural information processing systems (NeurIPS)*, 2017.<br><br>**Paper Discussion 1b.**<br>Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift.<br>Sergey Ioffe, Christian Szegedy.<br>*Proceedings of the International Conference on Machine Learning (ICML)*, 2015. |
| Monday, February 5<br>**In Person**<br>Feb 4  Feb **5**  Feb 6  Feb 7  Feb 8  Feb 9  Feb 10 | **Lecture #4: Kernels and Dimensionality Reduction.**<br>[Slides] [Demo Notebook]<br>• The kernel trick<br>• Gram matrix versus feature extraction: systems tradeoffs<br>• Adaptive/data-dependent feature mappings<br>• Dimensionality reduction |

| | |
|---|---|
| Wednesday, February 7<br>**In Person**<br><br>Feb 4 5 6 **7** 8 9 10 | **Paper Discussion 2a.**<br>Palm: Scaling language modeling with pathways.<br>Aakanksha Chowdhery et al.<br>*Journal of Machine Learning Research (JMLR)*, 2023.<br><br>**Paper Discussion 2b.**<br>Language models are few-shot learners.<br>Tom Brown et al.<br>*In Advances in neural information processing systems (NeurIPS)*, 2020.<br><br>**Due: Review of paper 1a or 1b.** |
| Monday, February 12<br>**In Person**<br><br>Feb 11 **12** 13 14 15 16 17 | **Lecture #5: Adaptive Methods & Non-Convex Optimization.**<br>[Slides] [Demo Notebook]<br>&bull; Adaptive methods<br>&bull; AdaGrad<br>&bull; Adam<br>&bull; Non-convex optimization<br><br>**Due: Programming Assignment 1.** |
| Wednesday, February 14<br>**In Person**<br><br>Feb 11 12 13 **14** 15 16 17 | **Paper Discussion 3a.**<br>Random features for large-scale kernel machines.<br>Ali Rahimi and Benjamin Recht.<br>*In Advances in Neural Information Processing Systems (NeurIPS)*, 2007.<br><br>**Paper Discussion 3b.**<br>Feature Hashing for Large Scale Multitask Learning.<br>Kilian Weinberger, Anirban Dasgupta, Josh Attenberg, John Langford and Alex Smola.<br>*Proceedings of the International Conference on Machine Learning (ICML)*, 2009.<br><br>**Released: Programming Assignment 2.** |
| Monday, February 19<br>**Online Only**<br><br>Feb 18 **19** 20 21 22 23 24 | **Lecture #6: Hyperparameter Optimization.**<br>[Slides] [Demo Notebook]<br>&bull; Hyperparameter optimization<br>&bull; Assigning parameters from folklore<br>&bull; Random search over parameters |
| Wednesday, February 21<br>**In Person**<br><br>Feb 18 19 20 **21** 22 23 24 | **Paper Discussion 4a.**<br>Random shuffling beats sgd after finite epochs.<br>Jeff Haochen and Suvrit Sra.<br>*Proceedings of the International Conference on Machine Learning (ICML)*, 2019.<br><br>**Paper Discussion 4b.**<br>Adam: A method for stochastic optimization.<br>Diederik Kingma and Jimmy Ba.<br>*Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.<br><br>**Due: Review of paper 3a or 3b.** |
| Monday, February 26 | **February Break: No classes.** |
| Wednesday, February 28<br>**In Person**<br><br>Feb 25 26 27 **28** 29 Mar 1 2 | **Paper Discussion 5a.**<br>Random search for hyper-parameter optimization.<br>James Bergstra and Yoshua Bengio.<br>*Journal of Machine Learning Research (JMLR)*, 2012.<br><br>**Paper Discussion 5b.**<br>Scaling laws for neural language models.<br>Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei.<br>*arXiv preprint arXiv:2001.08361*, 2020. |
| Monday, March 4<br>**In Person**<br><br>Mar 3 **4** 5 6 7 8 9 | **Lecture #7: Parallelism.**<br>[Slides] [Demo Notebook]<br>&bull; Hardware trends that lead to parallelism<br>&bull; Sources of parallelism in hardware<br>&bull; Data parallelism<br>&bull; Extracting parallelism at different places in the computation<br>&bull; Simple parallelism on multicore<br><br>**Due: Programming Assignment 2.** |

| | |
|---|---|
| Wednesday, March 6<br>**In Person**<br><br>Mar 3 Mar 4 Mar 5 **Mar 6** Mar 7 Mar 8 Mar 9 | **Paper Discussion 6a.**<br>Map-reduce for machine learning on multicore.<br>Cheng-Tao Chu, Sang K Kim, Yi-An Lin, YuanYuan Yu, Gary Bradski, Andrew Y. Ng, and Kunle Olukotun<br>*In Advances in Neural Information Processing Systems (NeurIPS)*, 2007.<br><br>**Paper Discussion 6b.**<br>Hogwild: A lock-free approach to parallelizing stochastic gradient descent.<br>Feng Niu, Benjamin Recht, Christopher Re, and Stephen Wright.<br>*In Advances in Neural Information Processing Systems (NeurIPS)*, 2011. |
| Monday, March 11<br>**In Person**<br><br>Mar 10 **Mar 11** Mar 12 Mar 13 Mar 14 Mar 15 Mar 16 | **Lecture #8: Distributed Learning.**<br>[Slides]<br>• Learning on multiple machines<br>• SGD with all-reduce<br>• The parameter server<br>• Asynchronous parallelism on multiple machines<br>• Decentralized and local SGD<br>• Model and pipeline parallelism<br><br>**Due: Review of paper 5a or 5b.** |
| Wednesday, March 13<br>**In Person**<br><br>Mar 10 Mar 11 Mar 12 **Mar 13** Mar 14 Mar 15 Mar 16 | **Paper Discussion 7a.**<br>Flashattention: Fast and memory-efficient exact attention with io-awareness.<br>Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré.<br>*In Advances in Neural Information Processing Systems (NeurIPS)*, 2022.<br><br>**Paper Discussion 7b.**<br>A System for Massively Parallel Hyperparameter Tuning.<br>Liam Li et al.<br>*Proceedings of the 2nd Conference on Machine Learning and Systems (MLSys)*, 2020. |
| Monday, March 18<br>**In Person**<br><br>Mar 17 **Mar 18** Mar 19 Mar 20 Mar 21 Mar 22 Mar 23 | **Lecture #9: Low-Precision Arithmetic.**<br>[Slides]<br>• Memory<br>• Low-precision formats<br>• Floating-point machine epsilon<br>• Low-precision training<br>• Scan order<br><br>**Due: Review of paper 6a or 6b.**<br><br>**In-class project feedback activity.** |
| Wednesday, March 20<br>**In Person**<br><br>Mar 17 Mar 18 Mar 19 **Mar 20** Mar 21 Mar 22 Mar 23 | **Paper Discussion 8a.**<br>Large scale distributed deep networks.<br>Jeff Dean et al.<br>*In Advances in Neural Information Processing Systems (NeurIPS)*, 2012.<br><br>**Paper Discussion 8b.**<br>Towards federated learning at scale: System design.<br>Keith Bonawitz et al.<br>*In Proceedings of the 2nd MLSys Conference (MLSys)*, 2019. |
| Monday, March 25<br>**In Person**<br><br>Mar 24 **Mar 25** Mar 26 Mar 27 Mar 28 Mar 29 Mar 30 | **Lecture #10: Inference and Compression.**<br>[Demo Notebook]<br>• Efficient inference<br>• Metrics we care about when inferring<br>• Compression<br>• Fine-tuning<br>• Hardware for inference<br><br>**Due: Review of paper 7a or 7b.**<br><br>**Due: Final project proposals.** |
| Wednesday, March 27<br>**In Person**<br><br>Mar 24 Mar 25 Mar 26 **Mar 27** Mar 28 Mar 29 Mar 30 | **Paper Discussion 9a.**<br>Gpipe: Efficient training of giant neural networks using pipeline parallelism.<br>Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, HyoukJoong Lee, Jiquan Ngiam, Quoc V. Le, and Yonghui Wu.<br>*In Advances in Neural Information Processing Systems (NeurIPS)*, 2019. |

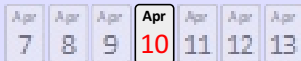| | |
|---|---|
| | **Paper Discussion 9b.**<br>Efficiently scaling transformer inference.<br>Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean.<br>*In Proceedings of Machine Learning and Systems (MLSys)*, 2023. |
| Monday, April 1 | **Spring Break: No classes.** |
| Wednesday, April 3 | **Spring Break: No classes.** |
| Monday, April 8<br>**In Person**<br>Apr 7 **Apr 8** 9 10 11 12 13 | **Lecture #11: Machine Learning Frameworks II.**<br>• Large scale numerical linear algebra<br>• Eager vs lazy<br>• ML frameworks in Python<br><br>**Due: Review of paper 8a or 8b.** |
| Wednesday, April 10<br>**In Person**<br>Apr 7 8 9 **Apr 10** 11 12 13 | **Paper Discussion 10a.**<br>Deep learning with limited numerical precision.<br>Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan.<br>*Proceedings of the International Conference on Machine Learning (ICML)*, 2015.<br><br>**Paper Discussion 10b.**<br>LoRA: Low-Rank Adaptation of Large Language Models.<br>Edward J. Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.<br>*Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. |
| Monday, April 15<br>**In Person**<br>Apr 14 **Apr 15** 16 17 18 19 20 | **Lecture #12: Hardware for Machine Learning.**<br>• CPUs vs GPUs<br>• What makes for good ML hardware?<br>• How can hardware help with ML?<br>• What does modern ML hardware look like?<br><br>**Due: Review of paper 9a or 9b.** |
| Wednesday, April 17<br>**In Person**<br>Apr 14 15 16 **Apr 17** 18 19 20 | **Paper Discussion 11a.**<br>Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding.<br>Song Han, Huizi Mao, and William J Dally.<br>*Proceedings of the International Conference on Learning Representations (ICLR)*, 2016.<br><br>**Paper Discussion 11b.**<br>GPTQ: Accurate post-training quantization for generative pre-trained transformers.<br>Frantar, Elias, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh.<br>*Proceedings of the International Conference on Learning Representations (ICLR)*, 2023. |
| Monday, April 22<br>**In Person**<br>Apr 21 **Apr 22** 23 24 25 26 27 | **Lecture #13: Modern Generative AI.**<br>• Scaling for large language models<br>• Challenges for LLM inference<br>• What does the future of generative AI look like?<br>• What are the policy and social implications of this technology?<br><br>**Due: Review of paper 10a or 10b.** |
| Wednesday, April 24<br>**Online Only**<br>Apr 21 22 23 **Apr 24** 25 26 27 | **Paper Discussion 12a.**<br>In-datacenter performance analysis of a tensor processing unit.<br>Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al.<br>*In Proceedings of the 44th Annual International Symposium on Computer Architecture (ISCA)*, 2017.<br><br>**Paper Discussion 12b.**<br>A Configurable Cloud-Scale DNN Processor for Real-Time AI.<br>Jeremy Fowers, Kalin Ovtcharov, Michael Papamichael, Todd Massengills, et al.<br>*In Proceedings of the 45th Annual International Symposium on Computer Architecture (ISCA)*, 2018. |
| Monday, April 29<br>**In Person**<br>Apr 28 **Apr 29** 30 May 1 May 2 May 3 May 4 | **Lecture #14: Large Scale ML on the Cloud.**<br>[Slides]<br>• Challenges of deployment<br>• Distributed learning at datacenter scale |

**Due: Review of paper 11a or 11b.**

**Due: Final project abstract draft. Can be submitted late until Wednesday afternooon; will discuss in class on Wednesday.**

| | |
|---|---|
| Wednesday, May 1<br>**In Person**<br><br>Apr 28 · Apr 29 · Apr 30 · **May 1** · May 2 · May 3 · May 4 | **Lecture #15: Final Project Disussion.** |
| Monday, May 6<br>**In Person**<br><br>May 5 · **May 6** · May 7 · May 8 · May 9 · May 10 · May 11 | **Lecture #16: Final Project Disussion.** |