

# CSCI 4701: Deep Learning

## Course Description

Deep Learning focuses on [artificial neural networks](#) and how they are trained and optimized. This course covers core concepts starting with backpropagation, including regularization and optimization techniques. Students will gain practical skills using the PyTorch framework and learn neural architectures used in both Computer Vision (CV) and Natural Language Processing (NLP), including Convolutional Neural Networks (CNNs) and Transformer-based models. The course introduces generative modeling with Variational Autoencoders (VAEs) and briefly covers current developments in deep learning, such as Latent Diffusion Models (LDMs) and Large Language Models (LLMs).

## Prerequisite(s)

- CSCI 4734: Machine Learning (official)
- Fundamental knowledge in Calculus, Linear Algebra, & Probability Theory
- Experience with Python & numpy (see this [tutorial](#))

## Contacts

Instructor(s)		
Name	Email	Office
Ismayil Shahaliyev	<a href="mailto:ishahaliyev@ada.edu.az">ishahaliyev@ada.edu.az</a>	B318

Academic Advisors <sup>1</sup>		
Name	Email	Program
Turkan Huseynova	<a href="mailto:thuseynova@ada.edu.az">thuseynova@ada.edu.az</a>	BSCS
Kamala Hashimova	<a href="mailto:khashimova@ada.edu.az">khashimova@ada.edu.az</a>	BSIT & BSMATH
Aysun Mustafazada	<a href="mailto:amustafazada@ada.edu.az">amustafazada@ada.edu.az</a>	BSCE & BSEEE

## Assessments

Item	Info	100%
Team Project	Students will apply the theoretical knowledge they gained to a practical semester-long team project.	25
Quizzes (8x)	Students will test their knowledge as they progress throughout the course.	25
Exam [Midterm]	Students will test their knowledge gained in the first half of the semester.	25
Exam [Final]	Students will test their knowledge gained throughout the semester.	25

## Technological Requirements

Students are expected to have [Github](#) and Google (for accessing [Google Colab](#)) accounts, and install the [git](#) version control system, [Python](#) (3.9 or later), and (optionally) [Jupyter](#). IDEs such as [VSCode](#) or [Cursor](#) are recommended for local development.

---

<sup>1</sup> Office: B building 1<sup>st</sup> floor. In case of **excused absences**, you are required to bring a document of proof to your academic advisor.

## Deadline Policy

Deadlines will never be extended. When submitting, a student must consider possible internet connection issues, uploaded file size, etc. **Late submissions** of one minute until thirty minutes will receive a 10% penalty, late submissions of thirty minutes until one day will receive a 25% penalty with no claim for bonuses (if there are any), late submissions of more than one day will receive zero. All deadlines are at 11 PM (23:00) Baku time (UTC+4), unless explicitly stated otherwise.

## Team Project

You will form teams of up to three members at the start of the semester. Students who are not part of a team one week before the project proposal deadline will be automatically assigned to a team at random. A **project proposal** that passes receives 10%. A proposal marked *revise* receives a provisional 0% until the required changes are submitted and approved, after which full credit is awarded. A proposal that is *rejected* and resubmitted as a substantially different proposal is eligible for a maximum of 5%. In addition to the proposal, there will be two further **project milestones**, scheduled before the midterm and final exam periods, each accounting for 45% of the total project grade. Detailed rubrics will be shared with you on blackboard.

## Quizzes

There will be **nine** computer-based quizzes assessing understanding of the material covered in the preceding lectures. The quiz grade will be calculated as the arithmetic mean of the **best eight** quizzes, with the lowest score automatically discarded (e.g. missing one quiz will not affect your grade). Each quiz will consist of **10 questions** with a **30-minute** time limit, including **4x** scenario-based multiple-choice questions (testing conceptual understanding and reasoning), **3x** short-answer questions (requiring precise explanations), and **3x** code-focused questions (e.g. identifying errors, predicting outputs, or completing small code fragments). One sample question per type will be shared with you before the first quiz.

## Exams

You will have two exams: midterm and final. Both exams will be computer-based (unless noted otherwise), consist of **20 questions** of the same types (**8x, 6x, 6x**) used in the quizzes, and have a **60-minute** time limit.<sup>2</sup> More specific rules and details will be communicated to you before the exam<sup>3</sup>.

## Suggested Reading

- Goodfellow, I., Bengio, Y., and Courville, A. Deep Learning. MIT Press, 2016.  
<https://www.deeplearningbook.org/>
- Zhang, A., Lipton, Z. C., Li, M., and Smola, A. J. Dive into Deep Learning.  
<https://d2l.ai>
- Study materials will be shared with you on blackboard throughout the semester.

## Course Outline

The content is subject to change. Please consistently check the course page on blackboard and the [ADA Academic Calendar](#) for modifications as well. The last day of the add/drop period, holidays, etc. are noted in the calendar.

<sup>2</sup> Quiz and exam dates, as well as project deadlines are noted in the [course outline](#) section.

<sup>3</sup> You are expected to be aware of the [Exam Rules and Regulations](#). If you have special needs or health issues, you are strongly recommended to contact the University's Student Academic Support Services well ahead of the examination date. It is your responsibility to manage conflicts in your schedule and notify your instructors about it at least two weeks before the examination date.

#	Topic	Learning Outcomes	Assessment
1	Deep Learning (DL) overview / Course structure	Describe the scope of DL and the course syllabus. Fulfill <b>technological requirements</b> .	
2	Mathematics of DL: Linear Algebra / Calculus	Work with vectors, matrices, and tensors, and perform basic <b>matrix operations</b> . Apply <b>norms</b> and inner products. Compute partial derivatives and apply the <b>chain rule</b> to multivariable functions. Optional: Understand intuition behind eigenvectors and SVD.	
3	Gradient Descent / Backpropagation I	Compute gradients efficiently using the chain rule on <b>computational graphs</b> . Perform forward and backward passes. Understand how <b>gradient descent</b> updates parameters to reduce a loss. Understand <b>automatic differentiation</b> of the PyTorch <a href="#">autograd</a> engine (with <a href="#">micrograd</a> ).	
4	Gradient Descent / Backpropagation II	<b>Implement</b> full backpropagation.	<b>Feb 3:</b> Quiz 1 [1-3]  Last day to submit team member details
5	Activation Functions / Neuron	Implement <b>activation functions</b> and understand <i>non-linearity</i> . Backpropagate over an <b>N-dimensional neuron</b> .	
6	Multilayer Perceptron (MLP)	Construct an MLP from stacked neurons. <b>Train a simple MLP classifier</b> on a small dataset.	<b>Feb 10:</b> Project proposal deadline
7	Images as Tensors / MLP on MNIST / Batching & Cross-Entropy	Understand <b>image representations</b> , tensor shapes, batching. Use <a href="#">torchvision datasets and dataloaders</a> . Train an MLP on MNIST with <b>SGD + Cross-Entropy</b> .	<b>Feb 12:</b> Quiz 2 [5-6]
8	Convolutional Neural Networks (CNN)	Define and implement 2D <b>cross-correlation</b> (convolution) and <b>pooling</b> with kernels, including padding and stride. <b>Train a LeNet-style CNN</b> on MNIST. Compare MLP with CNN.	
9	Mathematics of DL: Probability Theory	Describe <b>random variables</b> ; distinguish discrete and continuous distributions, work with <b>PMF/PDF</b> . Compute <b>expectation</b> , <b>variance</b> , and <b>covariance</b> . Understand <b>joint/marginal distributions</b> , <b>conditional probability</b> , <b>independence</b> , the <b>chain rule</b> , and <b>Bayes' rule</b> . Recognize the behavior of common <b>distributions</b> (e.g. Bernoulli, Multinoulli, Normal).	<b>Feb 19:</b> Quiz 3 [7-8]
10	Regularization	Apply <b>weight decay</b> and <b>dropout</b> to reduce overfitting. Handle <b>exploding</b> and <b>vanishing gradients</b> . Use <b>Xavier</b> and <b>He initialization</b> to improve gradient flow. Distinguish <b>local minima</b> from saddle points when reasoning about training dynamics.	
11	Optimization	Adjust <b>learning rate</b> and apply schedules. Use SGD with <b>momentum</b> to accelerate convergence. Apply <b>RMSProp</b> and <b>Adam</b> to handle noisy gradients. Compare optimizers based on convergence behavior and practical performance.	
12	Regularization / Optimization	<b>Train a regularized CNN on CIFAR-10</b> using optimizers. Apply <b>hyperparameter tuning</b> .	<b>Mar 3:</b> Quiz 4 [10-11]
13	Paper: AlexNet	Discuss AlexNet paper, its key ideas. Understand what is outdated. Discuss the paper structure.	
14	Bigram Model / Negative Log-Likelihood / Softmax	Build a character-level <b>bigram model</b> from text data and <b>sample</b> from it. Understand difference between probability and <b>likelihood</b> . Compute average <b>negative log-likelihood</b> as a loss. Explain the purpose of <b>softmax</b> .	
15	Neural Network N-gram Model / Mini-Batch Training	Construct a <b>neural N-gram model</b> . Train the model with mini-batch updates.	<b>Mar 12:</b> Quiz 5 [14]  Project milestone 1 deadline
16	<b>Midterm Exam [1-15] TUESDAY, MARCH 17</b>		
17	Midterm Exam Review	Half-semester overview.	
<b>HOLIDAYS, MARCH 20-30</b>			

18	Batch Normalization / Layer Normalization	Explain why normalization helps training deep networks. Implement <b>batch normalization</b> in a simple neural network and understand training vs evaluation behavior with <b>running statistics</b> . Understand the impact of <b>batch size</b> on batch normalization and when to prefer <b>layer normalization</b> instead.	
19	Residual Blocks / Residual Network for NLP	Understand the idea of <b>residual (skip) connections</b> and why they help very deep networks train. Modify a feed-forward N-gram model to include a <b>residual block</b> with appropriate dimensions. Interpret residual blocks as extending the function class while still allowing identity mappings, and connect them to vanishing gradients and regularization.	
20	Sequence Modeling: Autoregressive Models and RNN/LSTM	Understand <b>autoregressive sequence modeling</b> beyond fixed context windows. Explain how RNNs maintain state across time. Identify limitations of RNNs and LSTM/GRU (sequential computation, difficulty scaling, long-range dependency).	Apr 7: Quiz 6 [18-19]
21	Attention Mechanism	Understand <b>attention</b> as weighted information selection over a sequence. Derive and explain <b>queries</b> , <b>keys</b> , and <b>values</b> at the tensor level. Implement <b>attention</b> explicitly using matrix operations and verify shape consistency and normalization behavior.	
22	Transformer Architecture / Self-Attention	Understand <b>self-attention</b> as attention within a single sequence. Explain how <b>Transformer blocks</b> are constructed from self-attention, feed-forward layers, residual connections, and layer normalization. Explain why Transformers scale.	Apr 14: Quiz 7 [20-21]
23	Transformer Architecture / Transformer Blocks	Assemble a complete <b>Transformer block</b> from self-attention and feed-forward sublayers. Trace the forward and backward signal flow through the block. Analyze training stability, gradient propagation, and sensitivity to initialization and learning rate.	
24	Paper Reading: Transformer, Vision Transformer, Swin-Transformer	Extract the core architectural ideas from the Transformer, Vision Transformer, and Swin Transformer papers. Compare how attention is applied to sequences versus images and how architectural constraints affect scalability and efficiency.	Apr 21: Quiz 8 [22-23]
25	Mathematics of DL: Information Theory and Probabilistic Modeling	Compute entropy, cross-entropy, and <b>Kullback-Leibler (KL) divergence</b> for discrete and continuous distributions. Derive cross-entropy loss from <b>maximum likelihood</b> principles. Interpret common deep learning losses as probabilistic objectives.	
26	Variational Autoencoders I	Introduce <b>latent-variable generative models</b> . Understand the idea of latent representations and probabilistic <b>encoders/decoders</b> . Explain the role of <b>approximate inference</b> and why variational methods are needed. Analyze how VAEs learn structured latent spaces and how they differ from deterministic <b>autoencoders</b> .	April 28: Project milestone 2 deadline
27	Variational Autoencoders II	Understand the VAE training objective ( <b>ELBO</b> optimization). Implement <b>VAE</b> . Interpret <b>reconstruction</b> and <b>regularization</b> terms and their trade-off.	
28	Diffusion Models / Score Matching	Formulate diffusion models via <b>forward noising</b> and learned <b>reverse denoising</b> processes. Derive the <b>training objective</b> and interpret it as denoising <b>score matching</b> . Explain sampling as iterative probabilistic inference.	May 5: Quiz 9 [25-27]
29	Foundation Models and Modern Trends	Explain large-scale pretraining and <b>transfer learning</b> as the basis of <b>foundation models</b> . Examine GPT, BERT, and CLIP as representative architectures for language and multimodal learning, and <b>latent diffusion models</b> (LDMs). Examine their scaling behavior and practical limitations.	
Final Exam [1-29] TUESDAY, MAY 12			

## Revisions

v0

- No revisions yet