

Robot Learning

▼ 🕒 Deep Reinforcement Learning

- ✅ CS231n

- 🕒 David Silver's lecture

🔗 [Teaching - David Silver](#)

<https://www.davidsilver.uk/teaching/>

- 🕒 DL Book

🔗 [GitHub - MingchaoZhu/DeepLearnin...](#)

<https://github.com/MingchaoZhu/DeepLearning>

- 🕒 RL Book

🔗 <http://incompleteideas.net/book/the-...>

<http://incompleteideas.net/book/the-book-2nd.html>

▪ 🕒 Imitation Learning

🔗 [Learning to Imitate | SAIL Blog](#)

▼ 🕒 Robot Learning Related Papers & Implementation

🔗 [You awesome title | Write an awesom...](#)

- tensorflow

A deeplearning framework developed by Google (as opposed to pyTorch from Facebook), which is used by most of deep learning algorithms from Google researchers.

So, it is necessary to familiarize myself with modules and functions that are commonly used in deep learning models.

- vision language model

🔗 [CLIP: Connecting text and images](#)

Especially large models, are powerful to help robots understand instructions or feedbacks from the environment.

(e.g. OpenAI's CLIP: Contrastive Language-Image Pre-training)

- tokenization/token compress technique

Compress tokens so as to enable real-time inference even though the backbone is a high capacity Transformer.

- **language conditioned architecture**

Fusion of multimodal inputs (Image and Language), I need to understand how such a conditioning layer works (e.g., AAAI 18-FiLM: Visual Reasoning with a General Conditional Layer)

- **basic policy training pipeline**

Common practice to map from state to action (as opposed to Robotics Transformers, which do not rely on such kind of manipulation policies)

- **experimental details and interweaving of SayCan, RT-1 and RT-2**

When look into the experimental details of these papers, I found they are interweaving and it takes some time to figure out and understand the underlying reasons for model selection and experimental design.

Good feature: each self-produced module is paired with a file for Unit Test (useful for debugging, testing and maintaining)

Opportunity: multimodal data fusion happens in the embedding space. Simple, but may shrink the subtle useful information from inputs. How about different ways of multimodal data fusion?