# Initial Project Proposal

## Structure

The deliverable will be a piece of open-source software.

## Topic

In 2017, Yoshua Bengio submitted an arXiv article named *The Consciousness Prior* [1]. He proposed the *Consciousness Prior Theory*, which his intuition about how the consciousness works and how the *Global Workspace Theory* could be implemented using recently developed *Deep Learning* technology. The article contains a preliminary plan about how all the pieces should come together to create a consciousness that can facilitate the learning process. However, to my best knowledge, there are still no existing experiments that can support his *Consciousness Prior Theory*.

Implementing the whole framework in one semester will be too ambitious. So I am going to implement a minimum embodied toy RL model to showcase part of his theory, more precisely, the abstraction part which produces the high-level abstraction $h$.

In Bengio's paper, he thinks the high-level abstraction $h$ is produced by some representation RNN function $F$, so that

$$h_t = F(x_t, h_{t-1})$$

where $x_t$ is the *observation* at time $t$, and $h_t$ is the *unconscious representation state* (or *high-level state*) at time $t$. I think this is very similar to the $RL^2$ paper [2], where $h$ is the hidden states of the RNN.
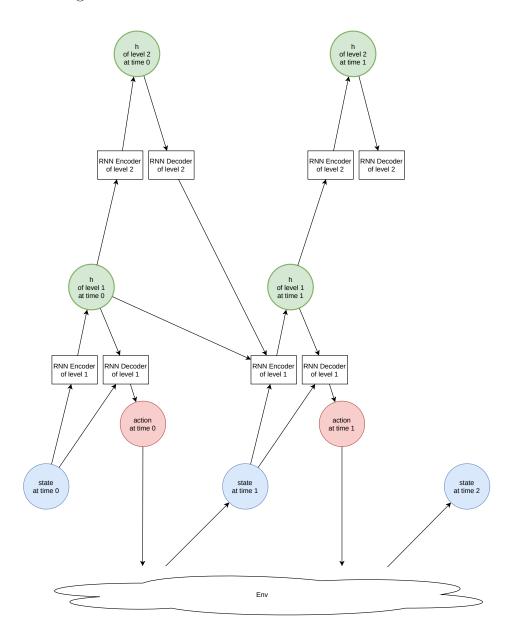
If we take the *Hierarchical RL* point of view, $h_t$ can be the *observation* of a high-level RL, and there could be multiple levels of $h$'s. (I am not sure how to automatically define *options* yet, maybe it can be another output of the function $F$.)

Moreover, the recent success of the Transformer model [3] indicates that it is possible to add Attention mechanism to extract $h$ more efficiently. I am going to combine Transformer and RNN if there's enough time.

I plan to use PyTorch to implement the neural networks, and I'll possibly use some other open-source projects (such as stable-baseline3) as sub-modules.

Hope it will work.

HRL diagram:

# References

[1] Yoshua Bengio. The Consciousness Prior. *arXiv:1709.08568 [cs, stat]*, December 2019. arXiv: 1709.08568.

[2] Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. RL\$ˆ2\$: Fast Reinforcement Learning via Slow Reinforcement Learning. *arXiv:1611.02779 [cs, stat]*, November 2016. arXiv: 1611.02779.

[3] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. *arXiv:1706.03762 [cs]*, December 2017. arXiv: 1706.03762.