# Grounded Language Learning in a Simulated 3D World

Spring 2021 CSE 6369 HLAI UTA

Ogenna Esimai

March 05, 2021

# Outline

- Disclaimer
- Topics
- Summary

# Outline

- **Disclaimer**
- Topics
- Summary

# Disclaimer

- Focusing on points

  - main

  - relevant

# Outline

- Disclaimer
- **Topics**
- Summary

# Topics (1)

- Need
  - communicate with agents

- Method
  - human language

- Scalable outcome
  - agent understanding of language
    - grounded

Source - Universal Intelligence: A Definition of Machine Intelligence. Legg et al. Minds & Machines (2007) 17:391–444

# Topics (2)

- Agent learns

  - interpret human language

  - simulated 3D environment

- "Grounded" caveat

# Topics (3)

- Agent–Environment Framework
  - Agent
    - a neural network
    - four inter-connected modules
      - convolutional **vision** module V
      - recurrent LSTM **language** module L
      - **mixing** module M
      - a two-layer LSTM **action** module A
  - Environment
    - DeepMind Lab, Beattie et al. (2016)
    - two connected rooms, each has two objects
    - other scenarios

# Topics (4)



Figure 3: **Unsupervised learning via auxiliary prediction objectives facilitates word learning**. Learning curves for a vocabulary acquisition task. The agent is situated in a single room faced with two objects and must select the object that correctly matches the textual instruction. A total of 59 different words were used as instructions during training, referring to either the shape, colours, relative size (larger, smaller), relative shade (lighter, darker) or surface pattern (striped, spotted, etc.) of the target object. **RP**: reward prediction, **VR**: value replay, **LP**: language prediction, **tAE**: temporal autoencoder. Data show mean and confidence bands (CB) across best five of 16 hyperparameter settings sampled at random from ranges specified in the appendix. Training episodes counts individual levels seen during training.

"Grounded Language Learning in a Simulated 3D World by Karl Moritz Hermann & Felix Hill et al."

9

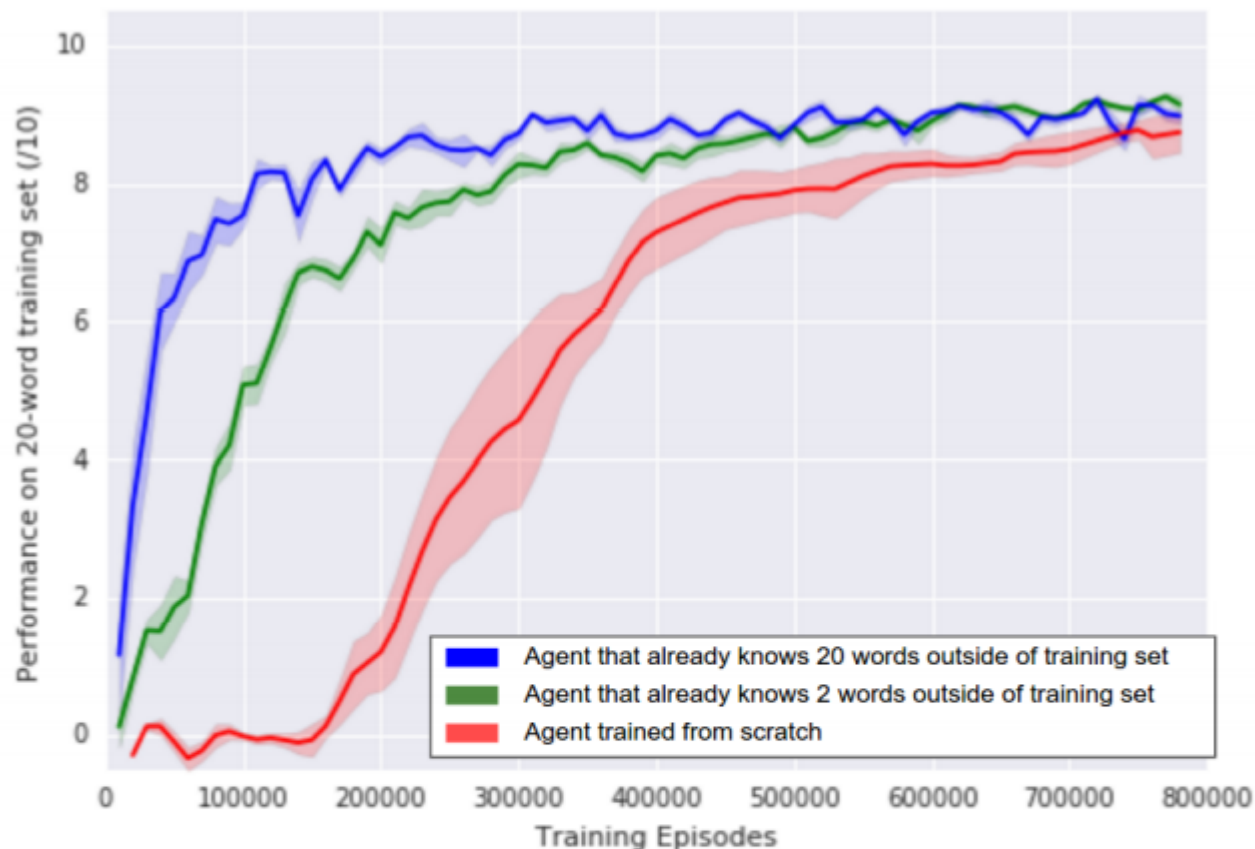# Topics (5)



Figure 4: **Word learning is much faster once some words are already known** The rate at which agents learned a vocabulary of 20 shape words was measured in agents in three conditions. In one condition, the agent had prior knowledge of 20 shapes and their names outside of the training data used here. In the second condition, the agent had prior knowledge of two shape words outside of the target vocabulary (same number of pre-training steps). In the third condition, the agent was trained from scratch. All agents used **RP**, **VR**, **LP** and **tAE** auxiliary objectives. Data show mean and confidence bands across best five of 16 hyperparameter settings in each condition, sampled at random from ranges specified in Appendix C.

# Topics (6)

- Extension of learning

  - curriculum

- Multi-task learning

  - Selection task - pick the X object or pick all X, where X denotes a colour term
  - Next to task - pick the X object next to the Y object, where X and Y refer to objects
  - In room task - pick the X in the Y room, where Y referred to the colour of floor in the target room

# Topics (7)



Figure 6: **Curriculum learning is necessary for solving more complex tasks.** For the agent to learn to retrieve an object in a particular room as instructed, a four-lesson training curriculum was required. Each lesson involved a more complex layout or a wider selection of objects and words, and was only solved by an agent that had successfully solved the previous lesson. The schematic layout and vocabulary scope for each lesson is shown above the training curves for that lesson. The initial (spawn) position of this agent varies randomly during training among the locations marked **x**, as do the position of the four possible objects among the positions marked with a white diamond. Data show mean and CB across best five of 16 randomly sampled hyperparameter settings in each condition.
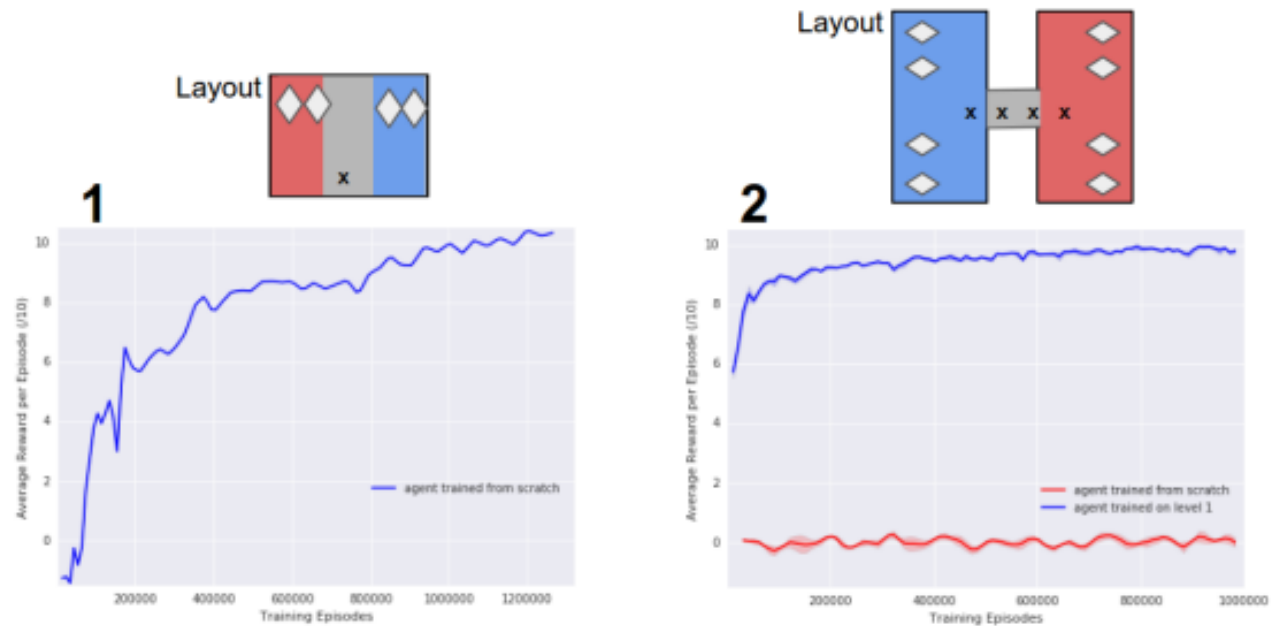
# Topics (8)



Figure 8: **Multi-task learning via an efficient curriculum of two steps.** A single agent can learn to solve a number of different tasks following a two-lesson training curriculum. The different tasks cannot be distinguished based on visual information alone, but require the agent to use the language input to identify the task in question.

Source: "Grounded Language Learning in a Simulated 3D World by Karl Moritz Hermann & Felix Hill et al."

# Topics (9)

- Final agent working

  - https://youtu.be/wJjdu1bPJ04

# Outline

- Disclaimer
- Topics
- Summary

# Summary

- Grounded learning in agent
  - role of unsupervised learning

- Concepts of language learning in agent
  - prior learning to compose vs. speed / bootstrap
  - generalization of knowledge
  - extension of learning / curriculum

- Future directions / Improvement
  - real world, non-simulated environment
  - does the order make a difference in performance? (Fig. 3)